

MULTISERVICE NETWORK PROCESSING DEVICE

Bidyut Parruck
Chulanur Ramakrishnan
Rami Zecharia
Nael Atallah
Hung Nguyen
Shankar Channabasappa

CROSS-REFERENCE TO RELATED APPLICATIONS

This application claims priority under 35 U.S.C. 120 from the following U.S. patent applications: 1) U.S. Patent Application Serial No. 09/528,802, filed March 20, 2000, 2) U.S. Patent Application Serial No. 09/539,479, filed March 30, 2000, 3) U.S. Patent Application Serial No. 09/539,306, filed March 30, 2000, 4) U.S. Patent Application Serial No. 09/539,478, filed March 30, 2000, 5) U.S. Patent Application Serial No. 09/539,461, filed March 30, 2000, 6) U.S. Patent Application Serial No. 09/539,476, filed March 30, 2000, and 7) U.S. Patent Application Serial No. 09/539,477, filed March 30, 2000. The subject matter of the above-listed seven patent applications is incorporated herein by reference.

BACKGROUND INFORMATION**The Maximus and Internet Growth**

As the Internet grows, the volume of data traversing the various networks that comprise the Internet has increased tremendously. While packet-based technology is widely believed to be the dominant technology going forward, in the short term, no single technology has emerged as the clear winner in the race to satisfy the seemingly insatiable demand for bandwidth, which is fueled by the spectacular growth of the Internet.

At the same time, carriers respond to customers' demand for new applications by rolling out increasingly sophisticated data services, some of which are highly delay-sensitive and/or bandwidth-intensive (e.g., voice or video). The introduction of these differentiated services opens up opportunities for competing system vendors to introduce new technologies and/or protocols. As a consequence, changes in the networks are characterized not only by vastly increased bandwidth requirements but also by a highly diverse array of technologies and protocols.

The aforementioned changes impact the **access, core/long-haul**, and **metro/regional** networks in different ways. In the foreseeable future, the demand for high speed **access** (i.e., in the megabit to low gigabit range) appears to be satisfied by a multitude of competing technologies, e.g., xDSL, cable modem, Frame Relay, T1, T3, ATM, and Gigabit Ethernet. Some of the newer technologies such as xDSL and cable modem have made impressive gains, enabling users to transmit and receive data at vastly increased data rates, thereby facilitating higher-bandwidth applications, some which were impractical only a few years ago when dial-up or lower speed ISDN connections were the dominant access technologies.

Corresponding changes also occur in the **core/long-haul** networks. In the core network, the emergence of an all-optical core and DWDM technologies have made it practical to transfer terabits of data across long-haul trunks at substantially lower costs.

1 However, all-optical technologies have not progressed to the point where it is possible to
2 process packets and/or cells in the optical domain. Processing intelligence is still very
3 much the province of electronics. Accordingly, an all-optical core, when operated in its
4 pure optical mode, is substantially devoid of processing intelligence, such as the type
5 required to correctly perform end-to-end routing of a packet or cell from one user to
6 another. Without resorting to the electrical plane, it is currently not practical to offer
7 differentiated data services (e.g., routing different streams of traffic in different ways in
8 accordance to their classes of service) via all-optical technologies.

9 The **metro/regional** network bears the full brunt of the powerful changing forces
10 that affect the access and core networks. As an aggregation mechanism for access
11 networks, the metro/regional network must move to higher speeds to cope with the
12 increased bandwidth requirements. The metro/regional network must also interoperate
13 with legacy and emerging transmission technologies and protocols, as well as adapt the
14 divergent technologies and protocols to its own transmission requirements and the
15 requirements of the core network. The metro/regional network must also perform
16 extensive cell and packet processing, as these tasks are no longer performed by the
17 optical core network. To facilitate differentiated service offerings, the metro/regional
18 network must also handle traffic management for different data streams having different
19 classes of service. In short, the metro/regional network of tomorrow is characterized by
20 higher speeds, diverse transmission technologies and protocols, and enhanced processing
21 requirements.

22 **Capabilities of the Maximus**

23 The Maximus is designed to capitalize on these trends. It is a high integration
24 ASIC designed for high speed switching and routing systems. The Maximus also finds
25 applications in access devices operating at OC-192, or as an aggregation device in the
26 core network. In a typical line card deployment, the Maximus is implemented between
27 an OC-192 or a quad OC-48 multi-service framer and the switch fabric.
28

29 **BRIEF DESCRIPTION OF THE DRAWINGS**

30 The present invention is illustrated by way of example, and not by way of
31 limitation, in the Figures of the accompanying drawings and in which like reference
32 numerals refer to similar elements and in which:
33

34 Figure 1 is a figure showing multi-service system.

35 Figure 2 is a figure showing OC-192 line card for traffic manager application.

36 Figure 3 is a figure showing OC-192 line card for traffic manager application.

37 Figure 4 is a figure showing various applications of the maximus chipset along with the
38 type field.

39 Figure 5 is a figure showing traffic manager applications.

40 Figure 5A is a figure showing uses of the Maximus chip with cell-based switch fabric.

41 Figure 6 is a figure showing SAR applications.

42 Figure 6A is a figure showing uses of the Maximus chip with packet-based switch fabric.

43 Figure 7 is a figure showing ingress tm (ATM=>ATM).

44 Figure 8 is a figure showing ingress tm (ATM=> MPLS packet).

45 Figure 9 is a figure showing ingress tm (MPLS packet=>ATM).

- 1 Figure 10 is a figure showing ingress tm (packet=>packet).
- 2 Figure 11 is a figure showing ATM encapsulation.
- 3 Figure 12 is a figure showing reassembly.
- 4 Figure 13 is a figure showing ingress packet bypass.
- 5 Figure 14 is a figure showing egress tm (ATM=>ATM).
- 6 Figure 15 is a figure showing egress tm (ATM=> MPLS packet).
- 7 Figure 16 is a figure showing egress tm (MPLS packet =>ATM).
- 8 Figure 17 is a figure showing egress tm (packet=>packet).
- 9 Figure 18 is a figure showing ATM de-encapsulation.
- 10 Figure 19 is a figure showing segmentation.
- 11 Figure 20 is a figure showing egress packet bypass.
- 12 Figure 21 is a figure showing Maximus chipset top level block diagram.
- 13 Figure 22 is a figure showing signals.
- 14 Figure 23 is a figure showing block addresses.
- 15 Figure 24 is a figure showing read access and timings definitions.
- 16 Figure 25 is a figure showing read access characteristic parameters.
- 17 Figure 26 is a figure showing write access and timings definitions.
- 18 Figure 27 is a figure showing write access characteristics parameters.
- 19 Figure 28 is a figure showing interface from the CPU if to the blocks.
- 20 Figure 29 is a figure showing global sync.
- 21 Figure 30 is a figure showing slot_count and CPU_port signals.
- 22 Figure 31 is a figure showing interface of the CPU if block.
- 23 Figure 32 is a figure showing internal implementation for locking the current CPU
- 24 access. parameters.
- 25 Figure 33 is a figure showing test mux block diagram.
- 26 Figure 34 is a figure showing soft resets.
- 27 Figure 35 is a figure showing input/output disable.
- 28 Figure 36 is a figure showing registers.
- 29 Figure 37 is a figure showing incoming spi-4 interface.
- 30 Figure 38 is a figure showing input control block in spii.
- 31 Figure 39 is a figure showing output control block in spii.
- 32 Figure 40 is a figure showing bit align block.
- 33 Figure 41 is a figure showing drop idle state machine.
- 34 Figure 42 is a figure showing control word format.
- 35 Figure 43 is a figure showing descriptions of fields in the control words.
- 36 Figure 44 is a figure showing control word type list.
- 37 Figure 45 is a figure showing rx state machine transition table.
- 38 Figure 46 is a figure showing extent of dip-4 codewords.
- 39 Figure 47 is a figure showing example of dip-4 encoding (odd parity).
- 40 Figure 48 is a figure showing example of a port calendar.
- 41 Figure 49 is a figure showing fifo status state machine.
- 42 Figure 50 is a figure showing example of dip-2 encoding (odd parity).
- 43 Figure 51 is a figure showing summary of internal memories.
- 44 Figure 52 is a figure showing interface ports.
- 45 Figure 53 is a figure showing register description.
- 46 Figure 54 is a figure showing outgoing spi-4 interface.

- 1 Figure 55 is a figure showing training sequence.
- 2 Figure 56 is a figure showing tx state machine transition table.
- 3 Figure 57 is a figure showing extent of dip-4 codewords.
- 4 Figure 58 is a figure showing example of dip-4 encoding (odd parity).
- 5 Figure 59 is a figure showing example of transfer data byte ordering.
- 6 Figure 60 is a figure showing tx state machine block.
- 7 Figure 61 is a figure showing port calendar.
- 8 Figure 62 is a figure showing status interface to database.
- 9 Figure 63 skew detection.
- 10 Figure 64 is a figure showing seg_table – 64 x 26 bits internal dual-port SRAM.
- 11 Figure 65 is a figure showing interface.
- 12 Figure 66 is a figure showing timing diagram-reassembly port empty.
- 13 Figure 67 is a figure showing timing diagram-back-to-back port transfer.
- 14 Figure 68 is a figure showing register description.
- 15 Figure 69 is a figure showing software interface registers.
- 16 Figure 70 is a figure showing read outgoing spi start-up parameters.
- 17 Figure 71 is a figure showing write outgoing spi start-up parameters.
- 18 Figure 72 is a figure showing read outgoing spi start-up parameters.
- 19 Figure 73 is a figure showing write outgoing spi start-up parameters.
- 20 Figure 74 is a figure showing read outgoing spi start-up parameters.
- 21 Figure 75 is a figure showing write outgoing spi start-up parameters.
- 22 Figure 76 is a figure showing write port parameters after initialization.
- 23 Figure 77 is a figure showing read port parameters after initialization.
- 24 Figure 78 is a figure showing read port full parameters.
- 25 Figure 79 is a figure showing write port full parameters.
- 26 Figure 80 is a figure showing read control status register.
- 27 Figure 81 is a figure showing write spio_tstmux_sel.
- 28 Figure 82 is a figure showing read spio_tstmux_sel.
- 29 Figure 83 is a figure showing lookup engine and neighboring blocks.
- 30 Figure 84 is a figure showing top level block diagram of the lookup engine.
- 31 Figure 85 is a figure showing data dependency on sop cell followed none sop cell on the
- 32 same port.
- 33 Figure 86 is a figure showing header selection.
- 34 Figure 87 is a figure showing label hashing & record retrieval.
- 35 Figure 88 is a figure showing port parameters, it is 64 entries.
- 36 Figure 89 is a figure showing encap table structure, it is 8 entries.
- 37 Figure 90 is a figure showing port type.
- 38 Figure 91 is a figure showing record memory, it is 2meg entries per memory.
- 39 Figure 92 is a figure showing parameter width.
- 40 Figure 93 is a figure showing external interface signals.
- 41 Figure 94 is a figure showing external interface signal names and direction.
- 42 Figure 95 is a figure showing global signals.
- 43 Figure 96 is a figure showing CPU interface to look up engine.
- 44 Figure 97 is a figure showing memory size.
- 45 Figure 98 is a figure showing memory size bits.
- 46 Figure 99 is a figure showing memory chip select.

- 1 Figure 100 is a figure showing memory chip select values.
- 2 Figure 101 is a figure showing address space.
- 3 Figure 102 is a figure showing registers addresses.
- 4 Figure 103 is a figure showing nop command.
- 5 Figure 104 is a figure showing read lut memory1 command.
- 6 Figure 105 is a figure showing write lut memory1 command.
- 7 Figure 106 is a figure showing read lut memory2 command.
- 8 Figure 107 is a figure showing write lut memory2 command.
- 9 Figure 108 is a figure showing read port memory command.
- 10 Figure 109 is a figure showing write port memory command.
- 11 Figure 110 is a figure showing setup connection command for key command.
- 12 Figure 111 is a figure showing teardown connection command for key command.
- 13 Figure 112 is a figure showing get hash FID command.
- 14 Figure 113 is a figure showing get hash FID response.
- 15 Figure 114 is a figure showing read cam memory1 command.
- 16 Figure 115 is a figure showing write cam memory1.
- 17 Figure 116 is a figure showing read cam memory2.
- 18 Figure 117 is a figure showing write cam memory2.
- 19 Figure 118 is a figure showing init all memories.
- 20 Figure 119 is a figure showing init memory 1 & cam 1.
- 21 Figure 120 is a figure showing init memory 2 & cam 2.
- 22 Figure 121 is a figure showing ATM types.
- 23 Figure 122 is a figure showing ATM MPLS tag location & format.
- 24 Figure 123 is a figure showing MPLS format.
- 25 Figure 124 is a figure showing lookup types.
- 26 Figure 125 is a figure showing header format.
- 27 Figure 126 is a figure showing data position within the header format.
- 28 Figure 127 is a figure showing header substitution.
- 29 Figure 128 is a figure showing special header structure.
- 30 Figure 129 is a figure showing encap table structure.
- 31 Figure 130 is a figure showing the encapsulation header replacement in l2 case.
- 32 Figure 131 is a figure showing hashing and memory accesses.
- 33 Figure 132 is a figure showing hash timing and lookup operations.
- 34 Figure 133 is a figure showing types of segmentation and applications for various traffic
- 35 types.
- 36 Figure 134 is a figure showing segmentation engine and neighboring blocks.
- 37 Figure 135 is a figure showing spi-4 data flow.
- 38 Figure 136 is a figure showing type 1 segmentation.
- 39 Figure 137 is a figure showing type 2 segmentation.
- 40 Figure 138 is a figure showing type 3 segmentation.
- 41 Figure 139 is a figure showing type 4 segmentation.
- 42 Figure 140 is a figure showing status cell access flow.
- 43 Figure 141 is a figure showing CPU access flow.
- 44 Figure 142 is a figure showing top-level block diagram for segmentation engine.
- 45 Figure 143 is a figure showing error conditions.
- 46 Figure 144 is a figure showing summary of internal memories.

- 1 Figure 145 is a figure showing seg_table – 65 x 39 bits internal dual-port SRAM
- 2 Figure 146 is a figure showing free_buffer – 80 x 7 bits internal SRAM.
- 3 Figure 147 is a figure showing q_fifo – 80 x 70 bits internal SRAM.
- 4 Figure 148 is a figure showing data SRAM – 640 x 64 bits internal SRAM.
- 5 Figure 149 is a figure showing crc_table – 65 x 32 bits internal dual-port SRAM.
- 6 Figure 150 is a figure showing statistic SRAM – 65 x 80 bits internal SRAM.
- 7 Figure 151 is a figure showing system.
- 8 Figure 152 is a figure showing segmentation ⇔ lookup.
- 9 Figure 153 is a figure showing segmentation ⇔ memory manager.
- 10 Figure 154 is a figure showing segmentation ⇔ PFQ.
- 11 Figure 155 is a figure showing segmentation ⇔ sch.
- 12 Figure 156 is a figure showing segmentation ⇔ CPU interface.
- 13 Figure 157 is a figure showing timing diagrams.
- 14 Figure 158 is a figure showing memory access.
- 15 Figure 159 is a figure showing registers.
- 16 Figure 160 is a figure showing summary of commands.
- 17 Figure 161 is a figure showing CPU inject packet.
- 18 Figure 162 is a figure showing read seg fbl memory.
- 19 Figure 163 is a figure showing write seg fbl memory.
- 20 Figure 164 is a figure showing read seg crc memory.
- 21 Figure 165 is a figure showing write seg crc memory.
- 22 Figure 166 is a figure showing read seg statistic memory.
- 23 Figure 167 is a figure showing write seg statistic memory.
- 24 Figure 168 is a figure showing read seg queue memory.
- 25 Figure 169 is a figure showing write seg queue memory.
- 26 Figure 170 is a figure showing read seg data memory.
- 27 Figure 171 is a figure showing write seg data memory.
- 28 Figure 172 is a figure showing read seg tbl memory.
- 29 Figure 173 is a figure showing write seg tbl memory.
- 30 Figure 174 is a figure showing memory manager block diagram.
- 31 Figure 175 is a figure showing internal memory manager block diagram.
- 32 Figure 176 is a figure showing multicast and BID/imBID release.
- 33 Figure 177 is a figure showing enqueue process in the internal memory manager.
- 34 Figure 178 is a figure showing dequeue process in the internal memory manager.
- 35 Figure 179 is a figure showing SDRAM manager block diagram.
- 36 Figure 180 is a figure showing memory controller configuration.
- 37 Figure 181 is a figure showing summary of memories.
- 38 Figure 182 is a figure showing summary of CAM.
- 39 Figure 183 is a figure showing segmentation engine.
- 40 Figure 184 is a figure showing per flow queue engine.
- 41 Figure 185 is a figure showing DDR-SDSRAM.
- 42 Figure 186 is a figure showing reassembly.
- 43 Figure 187 is a figure showing CPU interface.
- 44 Figure 188 is a figure showing global sync.
- 45 Figure 189 is a figure showing memory size is a figure showing 64 mbytes.
- 46 Figure 190 is a figure showing memory size is a figure showing 256 mbytes.

- 1 Figure 191 is a figure showing register description.
- 2 Figure 192 is a figure showing CPU_wr_SDRAM.
- 3 Figure 193 is a figure showing CPU_rd_SDRAM.
- 4 Figure 194 is a figure showing CPU_wr_fiml.
- 5 Figure 195 is a figure showing CPU_rd_fiml.
- 6 Figure 196 is a figure showing CPU_wr_cam.
- 7 Figure 197 is a figure showing CPU_rd_cam.
- 8 Figure 198 is a figure showing CPU_init_mm.
- 9 Figure 199 is a figure showing inputs from segmentation.
- 10 Figure 200 is a figure showing interface with SDRAM bw optimizer is a figure showing
- 11 SWM requests to write to SDRAM.
- 12 Figure 201 is a figure showing interface to SDRAM bw optimizer is a figure showing
- 13 write to SDRAM.
- 14 Figure 202 is a figure showing interface with per flow queue is a figure showing receive
- 15 de-queue from PFQ.
- 16 Figure 203 is a figure showing interface with SDRAM bw optimizer is a figure showing
- 17 SRM requests to read from SDRAM.
- 18 Figure 204 is a figure showing interface with SDRAM bw optimizer is a figure showing
- 19 SRM read from SDRAM.
- 20 Figure 205 is a figure showing interface with reassembly.
- 21 Figure 206 is a figure showing applications for various traffic types.
- 22 Figure 207 is a figure showing reassembly engine and neighboring blocks.
- 23 Figure 208 is a figure showing type 1 reassembly (ingress memory cells => switch cells).
- 24 Figure 209 is a figure showing type 1 reassembly (aal5 memory cells => switch cells).
- 25 Figure 210 is a figure showing type 1 reassembly (aal5 like memory cells => switch
- 26 cells).
- 27 Figure 211 is a figure showing type 2 reassembly (52 byte ATM cells => packet).
- 28 Figure 212 is a figure showing type 3 reassembly (48 byte ATM cells => packet).
- 29 Figure 213 is a figure showing type 4 reassembly (switch cells => packet).
- 30 Figure 214 is a figure showing type 5 reassembly (52 byte ATM cells => ATM cells with
- 31 translation).
- 32 Figure 215 is a figure showing type 6 reassembly (48 byte ATM cells => packet).
- 33 Figure 216 is a figure showing type 7 reassembly (48 byte cells => AAL-5 cells).
- 34 Figure 217 is a figure showing type 8 reassembly (switch cells => packet).
- 35 Figure 218 is a figure showing top level block diagram for the reassembly engine.
- 36 Figure 219 is a figure showing free buffer list block diagram.
- 37 Figure 220 is a figure showing crc_gen block diagram.
- 38 Figure 221 is a figure showing ATM MPLS tag location & format.
- 39 Figure 222 is a figure showing MPLS format.
- 40 Figure 223 is a figure showing ATM types.
- 41 Figure 224 is a figure showing seg_table – 64 x 18 bits internal single-port SRAM.
- 42 Figure 225 is a figure showing output_control_memory – 128 x 18 bit internal SRAM.
- 43 Figure 226 is a figure showing output_port_queue – 64 x 18 bit internal SRAM.
- 44 Figure 227 is a figure showing free_buffer – 128 x 7 bit internal SRAM.
- 45 Figure 228 is a figure showing data_SRAM – 1024 x 64 bit internal dual port SRAM.
- 46 Figure 229 is a figure showing header_SRAM – 256 x 64 bit internal dual port SRAM.

- 1 Figure 230 is a figure showing CPU_data_SRAM – 36 x 64 bit internal dual port SRAM.
- 2 Figure 231 is a figure showing crc_table – 64 x 32 bits internal dual-port SRAM.
- 3 Figure 232 is a figure showing statistic SRAM – 64 x 112 bit internal SRAM.
- 4 Figure 233 is a figure showing look-up SRAM – 4m x 36 bit external SRAM.
- 5 Figure 234 is a figure showing flow header format.
- 6 Figure 235 is a figure showing control word formats.
- 7 Figure 236 is a figure showing Maximus header format.
- 8 Figure 237 is a figure showing ATM header format.
- 9 Figure 238 is a figure showing MPLS format.
- 10 Figure 239 is a figure showing aal5 trailer format.
- 11 Figure 240 is a figure showing system.
- 12 Figure 241 is a figure showing reassembly ⇔ memory manager.
- 13 Figure 242 is a figure showing reassembly ⇔ SPI.
- 14 Figure 243 is a figure showing reassembly ⇔ CPU.
- 15 Figure 244 is a figure showing internal memory manager to reassembly.
- 16 Figure 245 is a figure showing reassembly to spi out - port empty and back to back.
- 17 Figure 246 is a figure showing timing diagram-back-to-back port transfer.
- 18 Figure 247 is a figure showing memories.
- 19 Figure 248 is a figure showing error conditions.
- 20 Figure 249 is a figure showing registers.
- 21 Figure 250 is a figure showing commands.
- 22 Figure 251 is a figure showing read external ras memory.
- 23 Figure 252 is a figure showing write external ras memory.
- 24 Figure 253 is a figure showing flow control format.
- 25 Figure 254 is a figure showing read internal data memory.
- 26 Figure 255 is a figure showing write internal data memory.
- 27 Figure 256 is a figure showing read four internal memories (partial crc1, partial crc2,
28 reassembly table, output queue).
- 29 Figure 257 is a figure showing write four internal memories (partial crc1, partial crc2,
30 reassembly table, output queue).
- 31 Figure 258 is a figure showing read port statistics.
- 32 Figure 259 is a figure showing write port statistics.
- 33 Figure 260 is a figure showing read CPU data memory.
- 34 Figure 261 is a figure showing write CPU data memory.
- 35 Figure 262 is a figure showing read 2 internal memories (free buffer list (fbl), output
36 control).
- 37 Figure 263 is a figure showing write 2 internal memories (free buffer list (fbl), output
38 control).
- 39 Figure 264 is a figure showing read internal header data memory.
- 40 Figure 265 is a figure showing write internal header data memory.
- 41 Figure 266 is a figure showing get CPU cell command.
- 42 Figure 267 is a figure showing initialize free buffer list and output port queue memories
43 command.
- 44 Figure 268 is a figure showing summary of lookup engine block tags for various traffic
45 types.
- 46 Figure 269 is a figure showing ATM header format.

- 1 Figure 270 is a figure showing MPLS label format.
- 2 Figure 271 is a figure showing ethernet frame format (IEEE 802.3 format).
- 3 Figure 272 is a figure showing IP (internet protocol) header.
- 4 Figure 273 is a figure showing frame relay (fr) layer 2 frame format (shown with IP
- 5 payload).
- 6 Figure 274 is a figure showing PPP full frame format for unnumbered mode operation.
- 7 Figure 275 is a figure showing flow control entry format.
- 8 Figure 276 is a figure showing tm to fabric flow control summary.
- 9 Figure 277 is a figure showing fabric to tm control summary.
- 10 Figure 278 is a figure showing internal block diagram.
- 11 Figure 279 is a figure showing graphic representation of an enqueue and dequeue FID.
- 12 Figure 280 is a figure showing operation sequence table.
- 13 Figure 281 is a figure showing enqueue block diagram.
- 14 Figure 282 is a figure showing dequeue block diagram.
- 15 Figure 283 is a figure showing free buffer manager block diagram.
- 16 Figure 284 is a figure showing random number generator.
- 17 Figure 285 is a figure showing memory size.
- 18 Figure 286 is a figure showing memory chip select.
- 19 Figure 287 is a figure showing data dependency.
- 20 Figure 288 is a figure showing RED memory.
- 21 Figure 289 is a figure showing RED memory byte enable.
- 22 Figure 290 is a figure showing statistics counter location 0.
- 23 Figure 291 is a figure showing statistics counter location 1.
- 24 Figure 292 is a figure showing statistics counter byte enable.
- 25 Figure 293 is a figure showing dequeue memory.
- 26 Figure 294 is a figure showing dequeue memory byte enable.
- 27 Figure 295 is a figure showing FID enqueue location 0.
- 28 Figure 296 is a figure showing FID enqueue memory location 1.
- 29 Figure 297 is a figure showing FID enqueue memory byte enable.
- 30 Figure 298 is a figure showing BID external data structure.
- 31 Figure 299 is a figure showing BID external data structure byte enable.
- 32 Figure 300 is a figure showing dequeue memory.
- 33 Figure 301 is a figure showing dequeue memory byte enable.
- 34 Figure 302 is a figure showing multicasting FID enqueue memory location 0.
- 35 Figure 303 is a figure showing multicasting FID enqueue memory location 1.
- 36 Figure 304 is a figure showing multicasting FID enqueue memory byte enable.
- 37 Figure 305 is a figure showing multicasting FID root memory.
- 38 Figure 306 is a figure showing FID dequeue leaf tunneling internal data structure.
- 39 Figure 307 is a figure showing FID enqueue leaf tunneling internal data structure location
- 40 0.
- 41 Figure 308 is a figure showing FID enqueue leaf tunneling internal data structure location
- 42 1.
- 43 Figure 309 is a figure showing FID enqueue leaf tunneling internal data structure byte
- 44 enable.
- 45 Figure 310 is a figure showing tunneling FID root data structure.
- 46 Figure 311 is a figure showing interface signals.

- 1 Figure 312 is a figure showing CPU interface to PFQ.
- 2 Figure 313 is a figure showing segmentation interface diagram.
- 3 Figure 314 is a figure showing segmentation interface.
- 4 Figure 315 is a figure showing DBS & CBWFQ interface.
- 5 Figure 316 is a figure showing DBS interface.
- 6 Figure 317 is a figure showing memory and PFQ interface.
- 7 Figure 318 is a figure showing memory and PFQ interface signals.
- 8 Figure 319 is a figure showing memory manager commands to per flow queue.
- 9 Figure 320 is a figure showing per flow queue commands to memory manager.
- 10 Figure 321 is a figure showing mem to PFQ commands.
- 11 Figure 322 is a figure showing PFQ to mem commands.
- 12 Figure 323 is a figure showing memory manager and the per flow queue interface.
- 13 Figure 324 is a figure showing PFQ and memory manager PFQ command structure.
- 14 Figure 325 is a figure showing PFQ to segmentation timing diagram.
- 15 Figure 326 is a figure showing address space.
- 16 Figure 327 is a figure showing registers space.
- 17 Figure 328 is a figure showing nop command.
- 18 Figure 329 is a figure showing read FID enqueue memory command.
- 19 Figure 330 is a figure showing write FID enqueue memory command.
- 20 Figure 331 is a figure showing read FID dequeue memory command.
- 21 Figure 332 is a figure showing write FID dequeue memory command.
- 22 Figure 333 is a figure showing read BID memory command.
- 23 Figure 334 is a figure showing write BID memory command.
- 24 Figure 335 is a figure showing read stat memory command.
- 25 Figure 336 is a figure showing write stat memory command.
- 26 Figure 337 is a figure showing setup connection command.
- 27 Figure 338 is a figure showing tear-down connection command.
- 28 Figure 339 is a figure showing init FID memories.
- 29 Figure 340 is a figure showing init BID memory.
- 30 Figure 341 is a figure showing timing that includes the FID enqueue, dequeue, BID and
- 31 statistics memories.
- 32 Figure 342 is a figure showing enqueue memory.
- 33 Figure 343 is a figure showing RED memory.
- 34 Figure 345 is a figure showing dequeue memory.
- 35 Figure 346 is a figure showing BID memory.
- 36 Figure 347 is a figure showing statistics memory.
- 37 Figure 348 is a figure showing basic block diagram.
- 38 Figure 349 is a figure showing interface to/from shaper.
- 39 Figure 350 is a figure showing general signals.
- 40 Figure 351 is a figure showing status interface from output spi.
- 41 Figure 352 is a figure showing flowid memory.
- 42 Figure 353 is a figure showing pkt memory.
- 43 Figure 354 is a figure showing port calendar memory.
- 44 Figure 355 is a figure showing shp_strict memory.
- 45 Figure 356 is a figure showing shp_out_port memory.
- 46 Figure 357 is a figure showing register descriptions.

- 1 Figure 358 is a figure showing read external FID memory.
- 2 Figure 359 is a figure showing write external FID memory.
- 3 Figure 360 is a figure showing setup FID connection.
- 4 Figure 361 is a figure showing read external pkt memory.
- 5 Figure 362 is a figure showing write external pkt memory.
- 6 Figure 363 is a figure showing read port calendar memory.
- 7 Figure 364 is a figure showing write port calendar memory.
- 8 Figure 365 is a figure showing read shp strict memory.
- 9 Figure 366 is a figure showing write shp strict memory.
- 10 Figure 367 is a figure showing read shp out memory.
- 11 Figure 368 is a figure showing write shp memory.
- 12 Figure 369 is a figure showing init shp_out_port memory.
- 13 Figure 370 is a figure showing init shp_strict memory.
- 14 Figure 371 is a figure showing timing diagram.
- 15 Figure 372 is a figure showing shaper.
- 16 Figure 373 is a figure showing scheduler.
- 17 Figure 374 is a figure showing PFQ.
- 18 Figure 375 is a figure showing ingress/egress mode of operations.
- 19 Figure 376 is a figure showing interface to/from scheduler.
- 20 Figure 377 is a figure showing scheduler's interface.
- 21 Figure 378 is a figure showing signal description.
- 22 Figure 379 is a figure showing simple round-robin algorithm.
- 23 Figure 380 is a figure showing weighted-round-robin scheduling.
- 24 Figure 381 is a figure showing creating smaller bursts with factor.
- 25 Figure 382 is a figure showing strict priority.
- 26 Figure 383 is a figure showing links between all levels of scheduling.
- 27 Figure 384 is a figure showing links between memories during input phase (the ingress
- 28 mode).
- 29 Figure 385 is a figure showing links between memories during output phase (the egress
- 30 mode).
- 31 Figure 386 is a figure showing links between memories during output phase (the ingress
- 32 mode).
- 33 Figure 387 is a figure showing external memory – next flowid.
- 34 Figure 388 is a figure showing internal memory – qos parameters.
- 35 Figure 389 is a figure showing internal memory – qos descriptors.
- 36 Figure 390 is a figure showing internal memory – port parameters.
- 37 Figure 391 is a figure showing internal memory – port descriptor.
- 38 Figure 392 is a figure showing internal memory – next port.
- 39 Figure 393 is a figure showing registers.
- 40 Figure 394 is a figure showing read FID next memory.
- 41 Figure 395 is a figure showing write FID next memory.
- 42 Figure 396 is a figure showing read qos parameters.
- 43 Figure 397 is a figure showing write qos parameters.
- 44 Figure 398 is a figure showing read qos descriptor.
- 45 Figure 399 is a figure showing write qos descriptor.
- 46 Figure 400 is a figure showing read port parameters.

- 1 Figure 401 is a figure showing write port parameters.
- 2 Figure 402 is a figure showing read port descriptor.
- 3 Figure 403 is a figure showing write port descriptor.
- 4 Figure 404 is a figure showing read port next.
- 5 Figure 405 is a figure showing write port next.
- 6 Figure 406 is a figure showing init qos.
- 7 Figure 407 is a figure showing init port.
- 8 Figure 408 is a figure showing global timing for input/output stages.
- 9 Figure 409 is a figure showing the input phase pipelines.
- 10 Figure 410 is a figure showing the output phase pipelines.
- 11 Figure 411 is a figure showing shaper/meter block interface.
- 12 Figure 412 is a figure showing traffic shaper/meter interface for dual-leaky-bucket.
- 13 Figure 413 is a figure showing traffic shaper/meter interface for single-leaky-bucket.
- 14 Figure 414 is a figure showing general signals.
- 15 Figure 415 is a figure showing progress of a leaky bucket algorithm.
- 16 Figure 416 is a figure showing links in the timing wheel.
- 17 Figure 417 is a figure showing example for shaping FID in the timing wheel.
- 18 Figure 418 is a figure showing calculating future time slot.
- 19 Figure 419 is a figure showing the usage of peak and sustained rates.
- 20 Figure 420 is a figure showing Rin and K factor.
- 21 Figure 421 is a figure showing credit calculation diagram.
- 22 Figure 422 is a figure showing calculation of threshold meter.
- 23 Figure 423 is a figure showing memory description.
- 24 Figure 424 is a figure showing shaper rate ID for shaper mode.
- 25 Figure 425 is a figure showing meter threshold for meter mode.
- 26 Figure 426 is a figure showing shaper slot memory.
- 27 Figure 427 is a figure showing FID mem1 for shaper.
- 28 Figure 428 is a figure showing FID mem1 for meter.
- 29 Figure 429 is a figure showing FID mem 2 for shaper.
- 30 Figure 430 is a figure showing FID mem2 for meter.
- 31 Figure 431 is a figure showing register description.
- 32 Figure 432 is a figure showing read FID1 memory.
- 33 Figure 433 is a figure showing write FID1 memory when using shaper.
- 34 Figure 434 is a figure showing read FID2 memory.
- 35 Figure 435 is a figure showing write FID2 memory.
- 36 Figure 436 is a figure showing read shp_slot memory.
- 37 Figure 437 is a figure showing write shp_slot memory.
- 38 Figure 438 is a figure showing read rate_id memory.
- 39 Figure 439 is a figure showing write rate_id memory.
- 40 Figure 440 is a figure showing setup connection command.
- 41 Figure 441 is a figure showing read FID1 memory.
- 42 Figure 442 is a figure showing write FID1 memory.
- 43 Figure 443 is a figure showing read FID2 memory.
- 44 Figure 444 is a figure showing write FID2 memory.
- 45 Figure 445 is a figure showing read threshold (internal) memory.
- 46 Figure 446 is a figure showing write threshold (internal) memory.

- 1 Figure 447 is a figure showing read modwrite meter.
- 2 Figure 448 is a figure showing setup connection command.
- 3 Figure 449 is a figure showing timing diagrams for meter.
- 4 Figure 450 is a figure showing block addresses.
- 5 Figure 451 is a figure showing read access and timings definitions.
- 6 Figure 452 is a figure showing read access timing and definitions table.
- 7 Figure 453 is a figure showing write access and timings definitions.
- 8 Figure 454 is a figure showing write access and timings definition table.
- 9 Figure 455 is a figure showing interface from the CPU IF to the blocks.
- 10 Figure 456 is a figure showing configuration modes.
- 11 Figure 457 is a figure showing global sync.
- 12 Figure 458 is a figure showing slot_count and CPU_port signals.
- 13 Figure 459 is a figure showing multiplexing the test busses from the blocks.
- 14 Figure 460 is a figure showing interface of the CPU IF block.
- 15 Figure 461 is a figure showing internal implementation for locking the current CPU
- 16 access parameters.
- 17 Figure 462 is a figure showing register description.
- 18 Figure 463 is a figure showing CPU interface signals.
- 19 Figure 464 is a figure showing block addresses.
- 20 Figure 465 is a figure showing CPUif block register.
- 21 Figure 466 is a figure showing PFQ block register.
- 22 Figure 467 is a figure showing read FID enqueue memory (72 bits) - up to 2m locations.
- 23 Figure 468 is a figure showing write FID enqueue memory (72 bits) - up to 2m locations.
- 24 Figure 469 is a figure showing read FID dequeue memory (36 bits)- up to 1m locations.
- 25 Figure 470 is a figure showing write FID dequeue memory (36 bits) - up to 1m
- 26 locations.
- 27 Figure 471 is a figure showing read BID dequeue memory (36 bits) - up to 16m
- 28 locations.
- 29 Figure 472 is a figure showing write BID memory (36 bits) - up to 16m locations.
- 30 Figure 473 is a figure showing read stat memory (72 bits) - up to 2m locations.
- 31 Figure 474 is a figure showing write stat memory (72 bits) - up to 2m locations.
- 32 Figure 475 is a figure showing setup connection command for FID.
- 33 Figure 476 is a figure showing tear-down connection command for FID.
- 34 Figure 477 is a figure showing init FID memories.
- 35 Figure 478 is a figure showing init BID memory.
- 36 Figure 479 is a figure showing database block registers.
- 37 Figure 480 is a figure showing read external FID memory.
- 38 Figure 481 is a figure showing write external FID memory.
- 39 Figure 482 is a figure showing setup FID connection.
- 40 Figure 483 is a figure showing read external pkt memory.
- 41 Figure 484 is a figure showing write external pkt memory.
- 42 Figure 485 is a figure showing read port calendar memory.
- 43 Figure 486 is a figure showing write port calendar memory.
- 44 Figure 487 is a figure showing read shp strict memory.
- 45 Figure 488 is a figure showing write shp strict memory.
- 46 Figure 489 is a figure showing read shp out memory.

1 Figure 490 is a figure showing write shp_out memory.
2 Figure 491 is a figure showing init shp_out_port memory.
3 Figure 492 is a figure showing init shp_strict memory.
4 Figure 493 is a figure showing shaper/meter block registers.
5 Figure 494 is a figure showing read FID1 memory.
6 Figure 495 is a figure showing write FID1 memory when using shaper.
7 Figure 496 is a figure showing read FID2 memory.
8 Figure 497 is a figure showing write FID2 memory.
9 Figure 498 is a figure showing read shp_slot memory.
10 Figure 499 is a figure showing write shp_slot memory.
11 Figure 500 is a figure showing read rate_id memory.
12 Figure 501 is a figure showing write rate_id memory.
13 Figure 502 is a figure showing setup connection command.
14 Figure 503 is a figure showing read FID1 memory.
15 Figure 504 is a figure showing write FID1 memory.
16 Figure 505 is a figure showing read FID2 memory.
17 Figure 506 is a figure showing write FID2 memory.
18 Figure 507 is a figure showing read threshold (internal) memory.
19 Figure 508 is a figure showing write threshold (internal) memory.
20 Figure 509 is a figure showing readmod write meter.
21 Figure 510 is a figure showing setup connection command.
22 Figure 511 is a figure showing sceduler block registers.
23 Figure 512 is a figure showing read FID next memory.
24 Figure 513 is a figure showing write FID next memory.
25 Figure 514 is a figure showing read qos parameters.
26 Figure 515 is a figure showing write qos parameters.
27 Figure 516 is a figure showing read qos descriptor.
28 Figure 517 is a figure showing write qos descriptor.
29 Figure 518 is a figure showing read port parameters.
30 Figure 519 is a figure showing write port parameters.
31 Figure 520 is a figure showing read port descriptor.
32 Figure 521 is a figure showing write port descriptor.
33 Figure 522 is a figure showing read port next.
34 Figure 523 is a figure showing write port next.
35 Figure 524 is a figure showing init qos.
36 Figure 525 is a figure showing init port.
37 Figure 526 is a figure showing parameters per flowid.
38 Figure 527 is a figure showing mx2 external memory.
39

40 **DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENTS**

41 The present invention will now be described in detail with reference to a few
42 preferred embodiments thereof as illustrated in the accompanying drawings. In the
43 following description, numerous specific details are set forth in order to provide a
44 thorough understanding of the present invention. It will be apparent, however, to one
45 skilled in the art, that the present invention may be practiced without some or all of these

specific details. In other instances, well known process steps and/or structures have not been described in detail in order to not unnecessarily obscure the present invention.

System Overview

The Maximus implements a multi-service approach with dynamic bandwidth allocation among the various flows carrying the different types of traffic. The Maximus device is designed such that it can be used in different applications that fit the needs of various system architectures. Two of the more popular applications are traffic manager and SAR solutions for OC-192 line cards in routing and switching systems. In one embodiment, the Maximus device is implemented in the form of two integrated circuit chips, designated the MX-1 and MX-2. In another embodiment, the entirety of the Maximus device is implemented on one integrated circuit chip (the circuitry of the MX-1 and MX-2 integrated circuits are integrated onto a single integrated circuit chip).

Figure 1 depicts an exemplary system level overview, illustrating how the Maximus device may enable OC-192 line cards in a router or switch system to receive different types of traffic, to dynamically allocate bandwidth among the various flows, to adapt the incoming traffic to the requirements of the cell-based or packet-based switch fabric, and to convert among the different traffic types for output. The term "router" as it is used in this document includes, but is not limited to, packet routers (for example, IP packet routers) and cell switches (for example, what are commonly called ATM switches).

In Figure 1, each line card can receive both ATM and packet traffic coming in on different flows across a fiber optic connection. The Maximus device on each line card, acting as an ingress device interfacing with a framer device, can receive, buffer and adapt the incoming traffic to the switching requirements of either the cell or packet switch fabric. On the egress side, the traffic can be sent out to the framer device either in the same traffic type (e.g., ATM in, ATM out) or in a different traffic type (e.g., ATM in, packet out or packet in, ATM out). With regard to dynamic bandwidth allocation, the Maximus offers extensive QoS capabilities to shape and schedule the flows as desired according to their QoS priorities.

Figures 2 and 3 below are two exemplary linecard views, illustrating how the Maximus device may be employed on an OC-192 line card operating with either a cell-based switch fabric (Figure 2) or a packet-based switch fabric (Figure 3). In these Figures, the traffic management or SAR solution is accomplished using two separate chips.

In Figure 2, a traffic management solution in conjunction with a cell-based switch fabric is illustrated. Four OC-48 optics couple with a quad OC-48 Framer to achieve an aggregate OC-192 data rate. This is not a requirement, however, as the Maximus Ingress Traffic Manager and the Maximus Egress Traffic Manager can also directly couple with an OC-192 framer to receive data at 10 Gbps. The Maximus supports 64 incoming or outgoing ports and up to 1 million separate flows in each direction. The interface between the Framer and Maximus Traffic Manager chips is SPI-4 (System Packet Interface-4), phase 2 (i.e., LVDS running at 400 MHz). Data storage and buffering employs SDRAMs, while pointers and tables are stored in SSRAMs. The Maximus Traffic Manager chips interface with a cell-based switch fabric in the example of Figure 2.

1 In the example of Figure 3, the Maximus Ingress Traffic Manager and the
2 Maximus Egress Traffic Manager are employed in the SAR configuration. For example,
3 the Maximus may receive ATM traffic, reassembles the ATM cells into packets along the
4 ingress path to switch through a packet-based switch fabric. On the egress path, the
5 packets may be segmented back into ATM cells for output to the quad OC-48 framers.
6 However, due to the multi-service capabilities offered, packets may also be received and
7 switched through the packet-based switch fabric. The data in these packets (whether
8 received as ATM cells or packets) may be output to the Framer as either ATM cells or
9 packets.

10 11 Applications

12 Figure 4 shows the various applications, along with their combinations of
13 incoming and outgoing traffic types.

14 Application types 0-3 in Figure 4 represent the ingress traffic manager
15 applications. Application types 8-11, on the other hand, represent the egress traffic
16 manager applications. These ingress and egress traffic management application types (0-
17 3 and 8-11) are graphically represented in Figure 5. Note that traffic management is
18 offered not only for packets but also for ATM cells and AAL5 cells in a single pair of
19 Maximus chips (ingress/egress pair). This reflects the multi-service approach of the
20 Maximus.

21 Figure 5A shows uses of the Maximus integrated circuit chip with a cell-based
22 switch fabric. The dashed box labeled "MX-1" indicates the contents of one integrated
23 circuit. Certain other parts of the integrated circuit (for example, the shaper block, the
24 scheduler block, and the memory management block) are omitted from Figure 5A in
25 order to clarify the illustration. These other parts are described in further detail below in
26 the later part of this patent document. The darkening of a block in Figure 5A indicates
27 that the block is substantially disabled in that it is not performing its primary function.
28 Blocks that are not darkened in Figure 5A are enabled to perform their primary functions.

29 Application types 4-6 and 12-14 represent the SAR application types, along with
30 their combinations of incoming and outgoing traffic types. These SAR application types
31 (4-6 and 12-14) are graphically represented in Figure 6. Note that these SAR application
32 types can handle not only segmentation of packets and reassembly of cells but can also
33 support ATM cell encapsulation and packet bypass (application types 4, 6, 12, and 14).
34 Again, this reflects the intelligent multi-service approach of the Maximus.

35 Figure 6A shows uses of the Maximus integrated circuit chip with a packet-based
36 switch fabric. The dashed box labeled "MX-1" indicates the contents of one integrated
37 circuit. Certain other parts of the integrated circuit (for example, the shaper block, the
38 scheduler block, and the memory management block) are omitted from Figure 6A in
39 order to clarify the illustration. These other parts are described in further detail below in
40 the later part of this patent document. The darkening of a block in Figure 6A indicates
41 that the block is substantially disabled in that it is not performing its primary function.
42 Blocks that are not darkened in Figure 6A are enabled to perform their primary functions.

43 The data formats for the various application types set forth in Figures 5, 5A, 6 and
44 6A are discussed below, with reference to the data format drawings of Figures 7-20.
45

Data Formats

Figure 7 illustrates the data format for Ingress Application Type 0 (ATM => ATM). In this application, the incoming data via the incoming SPI interface is ATM cells (52 bytes, which includes 4 bytes of ATM header and 48 bytes of data). Four bytes of pad is added prior to sending the incoming data to the lookup block (since the Maximus internal data path is 64-bit or 8-byte wide). The segmentation block adds 8 bytes of pad to allow the memory manager to pass the 64-byte memory cells to SDRAMs for storage. The reassembly block adds either 8 or 16 bytes of header to the 64-byte memory cell prior to sending to the cell-based switch (via the outgoing SPI block). The reader is invited to review egress Application type 8 (Figure 14) for the analogous process on the egress side, if the data is to be output to the Framer by the egress Maximus chip as ATM cells. Egress application type 8 is discussed later herein.

Figure 8 illustrates the data format for Ingress Application Type 1 (ATM => MPLS packet). In this application, the incoming data via the incoming SPI interface is ATM cells (52 bytes, which includes 4 bytes of ATM header and 48 bytes of data). Four bytes of pad is added prior to sending the incoming data to the lookup block (since the Maximus internal data path is 64-bit or 8-byte wide). The segmentation block removes the 4-byte ATM header and adds 12 more bytes of pad to allow the memory manager to pass the 64-byte memory cells to SDRAM for storage. The reassembly block adds either 8 or 16 bytes of header to the 64-byte memory cell prior to sending to the cell-based switch (via the outgoing SPI block). The reader is invited to review egress Application type 9 (Figure 15) for the analogous process on the egress side, if the data is to be output to the Framer by the egress Maximus chip as MPLS packets. Egress application type 9 is discussed later herein.

Figure 9 illustrates the data format for Ingress Application Type 2 (MPLS packet => ATM). In this application, the incoming data via the incoming SPI interface is MPLS packet (Nx16 bytes burst). The packet data is passed onto the lookup engine. The segmentation block segments the packet data into 48-byte AAL5 cells and appends 16 bytes of pad to allow the memory manager to pass the 64-byte memory cells to SDRAM for storage. The reassembly block adds either 8 or 16 bytes of switch header to the 64-byte memory cell prior to sending to the cell-based switch (via the outgoing SPI block). The reader is invited to review Application type 10 (Figure 16) for the analogous process on the egress side, if the data is to be output to the Framer by the egress Maximus chip as ATM cells. Egress application type 10 is discussed later herein.

Figure 10 illustrates the data format for Ingress Application Type 3 (Packet => Packet). In this application, the incoming data via the incoming SPI interface is packet (Nx16 bytes burst). The packet data is passed onto the lookup engine. The segmentation block segments the packet data into AAL5-like 64-byte cells. In this segmentation, CRC-32 is generated and an AAL5-like trailer is added to the end of the last 64-byte cell. The entire 64-byte AAL5-like cell is passed onto the memory manager for storage in SDRAM. The reassembly block adds either 8 or 16 bytes of switch header to the 64-byte memory cell prior to sending to the cell-based switch (via the outgoing SPI block). The reader is invited to review Application type 11 (Figure 17) for the analogous process on the egress side, if the data is to be output to the Framer by the egress Maximus chip as packets. Egress application type 11 is discussed later herein.

Figure 11 illustrates the data format for Ingress Application Type 4 (ATM Encapsulation). In this application, the incoming data via the incoming SPI interface is ATM cells (52 bytes, which includes 4 bytes of ATM header and 48 bytes of data). Four bytes of pad is added prior to sending the incoming data to the lookup block (since the Maximus internal data path is 64-bit or 8-byte wide). The segmentation block adds 8 bytes of pad to allow the memory manager to pass the 64-byte cell to SDRAM for storage. The reassembly block removes 8 bytes of pad, adds either 8 or 16 bytes of switch header prior to sending the resultant 64-byte packet to the packet-based switch (via the outgoing SPI block). The reader is invited to review egress Application type 12 (Figure 18) for the analogous process on the egress side, if the data is to be output to the Framer by the egress Maximus chip as ATM cells. Egress application type 12 is discussed later herein.

Figure 12 illustrates the data format for Ingress Application Type 5 (Reassembly). In this application, the incoming data via the incoming SPI interface is ATM cells (52 bytes, which includes 4 bytes of ATM header and 48 bytes of data). Four bytes of pad is added prior to sending the incoming data to the lookup block (since the Maximus internal data path is 64-bit or 8-byte wide). The lookup block removes the 4-byte ATM header and passes the 48-byte ATM data to the segmentation block, wherein 16 bytes of pad are added to allow the memory manager to pass the 64-byte cells to SDRAM for storage. The reassembly block maps the 48-byte ATM cells stored in memory as 64-byte cell into packets by stripping 16 bytes of pad from each memory cell and adding a switch packet header before sending the resultant packet to the packet-based switch via the outgoing SPI. The reader is invited to review egress Application type 13 (Figure 19) for the analogous process on the egress side, if the data is to be output to the Framer by the egress Maximus chip as ATM cells. Egress application type 13 is discussed later herein.

Figure 13 illustrates the data format for Ingress Application Type 6 (Ingress Packet Bypass). In this application, the incoming data via the incoming SPI interface is a packet burst (Nx16 bytes). The packet data is passed to the lookup block, which subsequently passes the packet data to the segmentation block. The segmentation block segments the packet data into AAL5-like 64-byte cells. In this segmentation, CRC-32 is generated and an AAL5-like trailer is added to the end of the last 64-byte cell. The entire 64-byte AAL5-like cell is passed onto the memory manager for storage in SDRAM. The reassembly block maps the 64-byte AAL5-like memory cell into packets by removing the AAL5 encapsulation and may add 8 or 16 bytes of switch packet header. The reassembled packet, including the 8/16 bytes of switch packet header is passed on to the packet-based switch via the outgoing SPI. The reader is invited to review egress Application type 14 (Figure 20) for the analogous process on the egress side, if the data is to be output to the Framer by the egress Maximus chip as packets. Egress application type 14 is discussed later herein.

Figure 14 illustrates the data format for Egress Application Type 8 (ATM =>ATM). In this application, the incoming data (from the switch) via the incoming SPI interface is a switch cell which comprises 8/16 bytes of switch header, 4 bytes of ATM header, 48 bytes of ATM data, and 12 bytes of pad. The switch header is removed by the lookup block, resulting in a 64-byte cell. The 64-byte cell is then passed through the segmentation block to the memory manager to be stored in SDRAM as a 64-byte memory cell. The reassembly block translates the ATM header, removes 8 bytes of pad (leaving 4

1 bytes of pad behind to facilitate sending the resultant 52-byte ATM cell + 4 bytes of pad
2 to the outgoing SPI interface (since the Maximus internal data path is 64-bit or 8-byte
3 wide). These 4 bytes of pad are later removed to allow the 52-byte ATM cell to be sent
4 out to the Framer (via the outgoing SPI block).

5 Figure 15 illustrates the data format for Egress Application Type 9 (ATM
6 =>MPLS packet). In this application, the incoming data (from the switch) via the
7 incoming SPI interface is a switch cell which comprises 8/16 bytes of switch header, 48
8 bytes of AAL5 data, and 16 bytes of pad. The switch header is removed by the lookup
9 block, resulting in a 64-byte cell. The 64-byte cell is then passed through the
10 segmentation block to the memory manager to be stored in SDRAM as a 64-byte memory
11 cell. The reassembly block strips the 16 bytes of pad, performs a MPLS operation of
12 adding the MPLS packet header based on the contents of the FlowID lookup table before
13 sending the resultant packet out to the Framer (via the outgoing SPI block).

14 Figure 16 illustrates the data format for the Egress Application Type 10 (MPLS
15 packets =>ATM). In this application, the incoming data (from the switch) via the
16 incoming SPI interface is a switch cell which comprises 8/16 bytes of switch header, 48
17 bytes of AAL5 data, and 16 bytes of pad. The switch header is removed by the lookup
18 block, resulting in a 64-byte cell. The 64-byte cell is then passed through the
19 segmentation block to the memory manager to be stored in SDRAM as a 64-byte memory
20 cell. The reassembly block adds the 4-byte ATM header and appends 4-byte of pad, as
21 well as strip the 16-byte of pad. These 4 newly added pad bytes are subsequently
22 removed to allow the 52-byte ATM cell to be sent to the Framer (via the outgoing SPI
23 block).

24 Figure 17 illustrates the data format for the Egress Application Type 11 (packet
25 => packet). In this application, the incoming data (from the switch) via the incoming
26 SPI interface is a switch cell comprising 64 bytes of AAL5-like data and an 8/16-byte
27 switch header. The switch header is removed by the lookup block, resulting in a 64-byte
28 cell. The 64-byte AAL5-like cell is then passed through the segmentation block to the
29 memory manager to be stored in SDRAM as a 64-byte memory cell. The reassembly
30 block maps the 64-byte AAL5-like cells into a packet, and adds a 8/16-byte packet
31 header. The reassembly block also checks the CRC in the AAL5 trailer before removing
32 the AAL5-like encapsulation. The resultant packet is subsequently passed on to the
33 Framer (via the outgoing SPI block).

34 Figure 18 illustrates the data format for the Egress Application Type 12 (ATM de-
35 encapsulation). In this application, the incoming data via the incoming SPI interface is a
36 switch packet, which comprises 8/16 bytes of switch header, 4 bytes of ATM header, 48
37 bytes of ATM data, and 4 bytes of pad. The switch header is removed by the lookup
38 block, resulting in 56-byte cell. The 56-byte cell is passed to the segmentation block,
39 which adds 8 bytes of pad to allow the memory manager to pass the 64-byte cell to
40 SDRAM for storage. The reassembly block translates the ATM header, strips out 8 bytes
41 of pad (leaving 4 bytes to allow passing through the 8-byte/64-bit data path from the
42 reassembly block to the outgoing SPI block). These last four bytes of pad will be
43 removed, allowing the outgoing SPI to send the 52-byte ATM cell to the Framer.

44 Figure 19 illustrates the data format for the Egress Application Type 13
45 (Segmentation). In this application, the incoming data via the incoming SPI interface is a
46 packet burst, including the 8/16 bytes of switch header. The packet data is passed onto

the segmentation block via the lookup block. At the segmentation block, the packet data is segmented into 48-byte AAL5 cells, and 16 bytes of pad are added to allow the memory manager to pass the 64-byte cells to SDRAM for storage. The reassembly block removes the 16 bytes of pad, adds the 4-byte ATM header to the 48-byte ATM data, and pads with 4 additional bytes to allow passing through the 8-byte/64-bit data path from the reassembly block to the outgoing SPI block. These last four bytes of pad will be removed subsequently, allowing the outgoing SPI to send the 52-byte AAL5 cell to the Framer.

Figure 20 illustrates the data format for the Egress Application Type 14 (Egress Packet Bypass). In this application, the incoming data via the incoming SPI interface is a packet burst, which comprises 8/16 bytes of switch header and 64 bytes of AAL5 data. The packet data is passed onto the segmentation block via the lookup block. At the segmentation block, the packet data is segmented into 64-byte AAL5-like cells for storage in SDRAM memory. The reassembly block maps the 64-byte AAL5-like cells into packet, and adds a packet header of 0-16 bytes. It also checks the CRC in the AAL5 trailer before removing the AAL5 encapsulation. The packet is then passed to the Framer via the outgoing SPI block.

Block Diagram

Figure 21 is a high-level block diagram of the Maximus 1 integrated circuit chip, the Maximus 2 integrated circuit chip, and various memory devices coupled to the Maximus 1 and 2 integrated circuit chips. All of these integrated circuits are present, in one embodiment, on a line card of a router. The functionality of the Maximus 1 (MX-1) and Maximus 2 (MX-2) integrated circuit chips may be disposed on two different integrated circuit chips as illustrated in Figure 21. Alternatively, the circuitry of the Maximus 1 and Maximus 2 integrated circuit chips is combined such that all that circuitry is disposed on a single integrated circuit chip.

Features

General Features

- Supports full duplex 10Gbps performance
- Protocols and internetworking
 - ATM
 - Operation and Management (OAM) cells
 - Signaling support
 - VP/VC translation
 - ATM to packet internetworking with Segmentation and reassembly
 - VC merging support
 - Packets
 - MPLS
 - MPLS tag switching
 - IP to ATM internetworking with AAL5 segmentation and re-assembly.
 - MPLS header generation. Used for MPLS tunneling and aggregation

- 1 • MPLS merging several flows from different ports merged
- 2 to single output.
- 3 ▪ Any packet with MPLS-like header
- 4 ○ Management support
- 5 ▪ SNMP, MIB counters
- 6 ▪ Alarms and Events
- 7 ○ Multicast support for both ATM and Packets
- 8 • Applications supported:
- 9 ○ Ingress Traffic Manager
- 10 ▪ ATM Cells in, ATM Cells out.
- 11 ▪ ATM (AAL5) Cells in, packets out.
- 12 ▪ MPLS Packets in, ATM (AAL5) Cells out.
- 13 ▪ MPLS Packets in, MPLS Packets out.
- 14 ○ Egress Traffic Manager
- 15 ▪ ATM Cells in, ATM Cells out.
- 16 ▪ ATM (AAL5) Cells in, Packets out.
- 17 ▪ MPLS Packets in, ATM (AAL5) Cells out.
- 18 ▪ MPLS Packets in, MPLS Packets out.
- 19 ○ Reassembler
- 20 ▪ ATM (AAL5) cells in, packets out
- 21 ▪ ATM cells in, ATM cell encapsulated inside packet out.
- 22 ▪ MPLS Packets in, MPLS Packets out.
- 23 ○ Segmentation
- 24 ▪ MPLS Packets in, ATM (AAL5) Cells out.
- 25 ▪ ATM cell encapsulated inside packet in, de-encapsulated ATM cell
- 26 out.
- 27 ▪ MPLS Packets in, MPLS Packets out.
- 28 • Interfaces supported
- 29 ○ Generic CPU Interface
- 30 ○ SPI-4 phase-2 interfaces for both packets and cells.
- 31 ▪ 16-bit, 400 MHz, DDR interface.
- 32 ▪ Supports up to 64 ports (any combination of STS-1 – STS-192).
- 33 ▪ SPI-4 used for both incoming and outgoing interfaces.
- 34 • Traffic management and QoS mechanisms
- 35 ○ Traffic Management on Ingress
- 36 ▪ Priority scheduling with weighted queuing
- 37 • Per Flow per egress port per class scheduling (up to 1k
- 38 queues)
- 39 • Flow control with the switch
- 40 ○ Traffic Management on Egress
- 41 ▪ Traffic shaping through dual leaky bucket
- 42 ▪ Traffic scheduling
- 43 • 64 Utopia level 4 ports
- 44 • 8 Class of service
- 45 • GCRA scheduling
- 46

- Weighted fair queuing
- VP tunneling support
- CBR tunneling
- MPLS tunneling

CPU Interface

- Supports CPU interface for management and control information.
- Supports 32-bit CPU bus interface.
- Support Mux/Non-Mux and Big Endian/Little Endian modes
- Supports interrupt generation.
 - Hierarchical interrupt structure

Incoming SPI-4 Interface

- Supports the SPI-4 phase II specifications for 10 Gbps rate.
- Unpacks incoming 16-bit DDR data at 400 MHz into 64-bit wide word at 200 MHz.
- Performs diagonal interleaved parity checking.
- Extracts in-band port address, SOP, EOP and error-control code.
- In-band start-of-FIFO status signal.
- 200 MHz for the FIFO status transfer.
- Supports maximum 64 ports.
- Generates 2-bit FIFO status back to the transmit device using programmed port calendar.
 - Programmable port calendar with STS-1 granularity.

Lookup Engine

- Supports up to 1M flows.
- Supports up to 64 input ports.
- Tables and configuration registers loaded via the processor interface.
- Lookup Algorithm
 - Label generator selects any 6 bytes from the first 32 bytes of the packet.
 - Label generator is programmable on a per port basis.
 - The label generator has one byte of granularity.
 - Hash generator creates 2 hash values or slots in parallel from a key using 2 separate CRC polynomials as hash functions.
 - Four bins per bucket to minimize collisions.
 - Hash table lookup constrained to resolve within 4 tries.
 - Overflow CAM of 256 or 64 entries, the largest available, for lookups that fail to resolve in 4 tries.
 - Fully pipelined, one lookup every 8 cycles.
- Supports the following flows with the following headers:
 - MPLS with PPP encapsulation.
 - MPLS with LLC/SNAP encapsulation.
 - MPLS with Frame Relay encapsulation..
 - Custom switch headers for traffic manager application.
 - Supports bypass lookup using MPLS-like header.

Segmentation Engine

- Segments incoming packet data into cells on a per port basis

- Supports a maximum of 64 external ports
- Supports an additional port for data from processor
- Segments incoming packet data into ATM cells
- Appends AAL-5 trailer for ATM segmentation
- Segments incoming packet data into Switch cells of 64 bytes
- Appends AAL-5 like trailer for non-ATM segmentation
- Supports programmable maximum packet length of 16K bytes.
- Gathers per port statistics
 - 32-bit counter for incoming packets.
 - 48-bit counter for byte count.
 - 16-bit counter to count discarded packets.
- Tags bad packets from SPI, to be discarded by per flow queue.
- Internal SRAM to store per port segmentation context and data

Memory Manager

- Uses two channels of 128-bit (64-bit DDR) SDRAM
- Manages read and write requests to all eight banks separately to maximize utilization.
- Uses 128kb of internal SRAM for temporary storage of data to and from SDRAM.
- Sophisticated load-sharing algorithm to distribute traffic over various banks.

Reassembly and Header Processing Engine

- Types of operation:
 - ATM cells to ATM cells no action required with translation
 - ATM cells to ATM cells with translation
 - ATM cells to Switch cells
 - Switch cells to Switch cells
 - Switch cells to Packet Data
 - ATM cells to Packet Data
 - AAL-5 cells to ATM cells, appending ATM header
 - AAL-5 cells to Packet Data

Header Processing Features

- Supports a maximum of 1M total flows.
- Types of lookup supported
 - No lookup, pass-through traffic
 - ATM processing
 - ATM Header append
 - VPI, VCI header translation
 - Encapsulate raw cells and OAM cells into packets for routing through packet fabrics.
 - Encapsulate ATM cells into packets for routing through packet fabric.
 - MPLS processing
 - MPLS push/pop
 - MPLS header translation

Reassembly Features

- Reassembles incoming cells into packets on a per port basis
 - Supports a maximum of 64 external ports
 - Supports an additional port for data to processor
 - Reassembles incoming ATM cells into packet data.
 - Strips AAL-5 trailer for ATM reassembly
 - Segments incoming Switch cells of 64 bytes into packet data.
 - Strips AAL-5 like trailer for non-ATM reassembly.
- Supports programmable maximum packet length of 16K bytes.
- Gathers per port statistics
 - 32-bit counter for outgoing packets.
 - 48-bit counter for byte count.
 - 16-bit counter to count discarded packets.
- Tags bad packets to SPI, to be discarded by downstream device.
- Internal SRAM to store per port segmentation context and data

Outgoing SPI-4 Interface

- Supports SPI-4 phase II spec at 10 Gbps rate.
- Packs the 64-bit data at 200 MHz to 16 bit DDR data at 400 MHz.
- Multiplexes the data and the control onto the SPI-4 bus.
- Diagonal Interleaved parity generation.
- Programmable data transfer burst sizes of multiple 16 bytes.
- Passes the FIFO status information to the scheduler to stop transferring.
- Schedules training sequence for data path de-skew.
- Supports maximum 64 ports
- Receives 2-bit FIFO status from the receive device using programmed port calendar.
 - Programmable port calendar with STS-1 granularity.

Per Flow Queue Engine

- See below, page 88.

Traffic Shaper and Traffic Meter

- See below, page 129.

Output Scheduler

- See below, page 118.

Signals

See Figure 22.

MX 1 CPU Interface

This block is responsible for three different functions:

- Connects between the blocks and the CPU via the CPU interface. It decodes the address to Figure which block to access and generate the chip select to each block.
- Generates global synchronization signals to the rest of the blocks.
- Holds general control registers like version register, mode register etc.

Features

1. Address bus of 10 bits per chip.
 Bits 0-5 are used as an address inside the block (up to 64 registers per block)
 Bits 9-6 are used to select a block (up to 16 blocks per chip)
2. Data bus of 32 bits per chip.
 Externally this bus is b-directional
 Internally this bus is split to 2 busses: DATA_IN and DATA_OUT.
3. Chip select signal is the qualifier for data, address and read/write busses.
4. The CPU interface can support multiplexed CPU interface or non-multiplexed CPU interface. When using multiplexed interface, the address and data are folded to the DATA[31:0] bus.
5. CPU accesses all registers in all blocks via direct read/write commands.
6. CPU accesses all internal/external memories in all blocks via indirect accesses using the general-purpose registers dedicated for these accesses.
7. Blocks addresses can be found in Figure 23.
8. Generates a global sync signal once every 520 cycles of the 200MHz clock (65 ports, 8 cycles per port).

Functional Description

The CPU IF block is used to interface with the CPU. The CPU IF block decodes the incoming address bus to generate a unique chip-select signal to each block. The CPU IF block also holds general registers that contain general parameters and modes for the entire chip. Basically any logic that is not related directly to a specific block is in this block.

The CPU IF block can operate in two modes depending on an input pins to select between them:

- Regular CPU access: separate busses for data and for address. In this case address bus is qualified by the chip-select signal.
- Multiplexed CPU access: one bus is multiplexed for address and for data. In this case address is qualified by the ALE signal.

The CPU interface is based on Intel's microprocessor interface (when multiplexing address/data bus). The following timing diagrams defines the characteristics of the CPU interface for Read and for Write accesses:

Read access characteristics

See Figures 24 and 25.

- RDWR_L signal has the same timings as the ADDR bus.
- A valid read cycle is defined when RDWR_L is latched as high on the negative edge of the CS_L signal.
- In non-multiplexed address/data bus architecture, ALE should be held high so parameters t_{Salr} , t_{Halr} , t_{VL} , t_{Slr} and t_{Hlr} are not applicable.
- In our implementation the address is latched on the negative edge of the CS_L therefore t_{Har} is not applicable.

Write access characteristics

See Figures 26 and 27.

- RDWR_L signal has the same timings as the ADDR bus.
- A valid write cycle is defined when RDWR_L is latched as low on the negative edge of the CS_L signal.
- In non-multiplexed address/data bus architecture, ALE should be held high so parameters t_{Salw} , t_{Halw} , t_{VL} , t_{Slw} and t_{Hlw} are not applicable.
- In our implementation the address is latched on the negative edge of the CS_L therefore t_{Haw} is not applicable.

Timing diagram of signals to blocks

The diagram in Figure 28 describes the timings of the CPU interface to each block. Note that chip-select signal is now asserted for only one cycle of the 200MHz clock. For a write or read cycle, all parameters are valid and stable before chip-select is asserted. When a read cycle occurs data is driven from the block after a maximum of 3 clock cycles of 200MHz clock.

GSYNC - Synchronization signal

GSYNC signal is a one-cycle pulse in the 200MHz clock, which is asserted once every 520 cycles. The SYNC signal is used for synchronization between the blocks and between the chips as well as to generate the slot-count and the CPU-port indication. See Figure 29.

Slot count and CPU port indication

The signals slot_count and CPU_port are generated inside each block base on the GSYNC signal. Each block first needs to sample GSYNC and use the sampled version to generate the signals. See Figure 30.

Interface

The CPU IF block generates chip-select signal per block, the rdwr_l, addr and data busses are propagated to the blocks directly. The blocks perform the read or write command based on the rdwr_l signal only when the chip-select signal per the block is asserted. The chip-select signal per each block is synchronous to the 200MHz clock domain already. Each block first samples the incoming signals and then uses them. This is being done to solve any timing issues that may occur. If during placement and routing the signals from the CPU IF to the blocks will violate the 5nSec budget, another set of FFs can be placed in the middle of the path.

For a read cycle the CPU IF block selects the data busses coming in from the blocks based on the chip-select per block (the address). See Figure 31.

The address, data and rd/wr signals takes about 4 cycles of 200MHz to propagate from the chip's pins to the block (6 cycles are used for delay calculation). 2 cycles are needed to perform read or write accesses in the block.

Therefore a total of 8 cycles are needed to perform a read or a write access.

Block Descriptions

CPU Interface

The negative edge of the CS_L signal qualifies the ADDR[9:0] and the RDWR_L signals. The positive edge of the CS_L signal qualifies the DATA bus for read and for write accesses. The CS_L is basically an enable signal to latches. The circuit shown in Figure 32 shows the basic latching circuit and the synchronization to 200MHz clock domain:

Test Mux

The Test Mux block is used to multiplex test mux outputs from all the blocks inside the chip. The output bus of the test mux go out to the pins that will be used for testing and debugging purposes. Refer to the Figure below for the Test Mux block diagram. The Tstmux_Blk_Sel register (refer to the Register section) is used to select which block signals to be output to the test pins (Tstmux_pin[31:0]). Each block inside the chip also has selection register (e.g. Seg_Tstmux_Sel, Ras_Tstmux_Sel) to select a group of internal signals to be observed.

The test mux outputs of all the blocks should be clocked out of flip-flops (refer to the diagram below). The CPU_clk that is running at 200 Mhz, is used to clock the output of the test mux inside the CPU Interface block. It is also going out to a test pin called Tstmux_clk. See Figure 33.

Soft Reset

Each block inside MX-1 chip can be reset individually by writing a '1' to its Soft_Reset register (refer to the Register section). Writing a '0' to the register will release the block from being reset and get it back to its normal operation. Each soft reset signal of each block is also coupled with the global power-on reset pin (POR_L). Refer to Figure 34 below.

Input/Output Disable

Each block in the MX-1 chip has 2 inputs coming from the CPU Interface block that are used to disable the input and output of the block. The Input_Disable and Output_Disable registers (refer to the Register section) can be programmed to disable the input or/and output of a particular block. Refer to Figure 35 below.

Registers

The following table describes the registers implemented in the CPU IF block. These registers act as control to the rest of the chip. The registers are implemented in the 200MHz clock domain even if there is no need to do so since their value is fixed and stable during normal operation. See Figure 36.

Incoming SPI-4

The Incoming System Packet Interface Level 4 (SPI-4) block receives 16-bit DDR ATM or Packet over SONET/SDH (POS) data from the Transmit Device (PHY or Switch based on ingress or egress application) and sends out-of-band FIFO status information to the Transmit Device. The incoming SPI-4 interface block outputs data

with control information to the lookup engine, and receives FIFO status from segmentation engine. Furthermore, the CPU can program the registers in the incoming SPI-4 interface block via CPU interface.

Features

- Supports the SPI-4 phases II spec at 10 Gbps rate
- Packs incoming 16-bit DDR data at 400MHz into 64-bit wide word at 200MHz
- Free buffer list to manage per port FIFO and absorb data while look up engine is doing the look up
- Diagonal interleaved parity checking for data
- Extracts in-band port address, SOP, EOP and error-control code and presents to the lookup engine
- Passes 2-bit FIFO status back to the transmit device in a round robin fashion based on the port calendar
- Programmable port calendar to output FIFO status
- In-band start-of-FIFO status signal
- 400 MHz for the FIFO status transfer
- De-skews +/- 1 bit time after running the training patterns
- Supports maximum 64 ports

Functional Description

The incoming SPI-4 interface block receives data with in-band port address, delineation information and error control code from an external transmit device. It separates the data and control and passes them to the lookup engine. It also sends out-of-band flow control information to the transmit device when it receives the information from the segmentation block. See Figures 37, 38 and 39.

Block Description

Bit Align

The bit-align block adds programmable delays to both data and control lines to correct for relative skew differences of up to +/- 1 bit time. The programmable delays are controlled by the bit align control block which receives pre-defined training patterns to determine the adjustment. It outputs 32 bit SDR data with two control signals. The control signals indicate which 16-bit word is data or control word. See Figure 40.

Bit Align Control

This block compares the received training patterns with the pre-defined patterns and determines the receive timing of the lines. Based on the timing, it generates the control signal the select different delays to de-skew the data and control lines.

Drop IDLE

This block drops the extra idle control words from the incoming data and writes only data, training pattern, and the first detected idle control word into the input FIFO. This block has a state machine which powers up in the idle state and remains in that state

until data or the first idle control word is detected. For the first detected idle control word, the state machine advances to the state named FIRST IDLE; and for non-idle control word it advances to the state named SCAN. If the state machine is in the FIRST IDLE state, it advances to the SCAN state if non-idle control word is detected, otherwise, it goes back to the IDLE state and asserts first_idle signal. If the state machine is in the SCAN state, it stays in this state until idle control word is detected and it then advances to the FIRST IDLE state. The incoming data and control words are allowed to pass through and written together with the control signals into the input FIFO when the state machine is in the SCAN or FIRST IDLE states. See Figure 41.

Input FIFO

This input FIFO is a 256x34 dual-port SRAM running at 400MHz to absorb incoming SPI-4 data, control words, and control signals. The write operation is controlled by the drop IDLE state machine, and the read operation is controlled by the receive state machine (RX_SM) when the FIFO is not empty.

Rx State Machine

The Rx state machine receives 32-bit words from the input FIFO with indications of data or control words of the two 16-bit words. It extracts the in-band control information, packs the data into 64-bit wide words, passes the training patterns to the bit align control block. It also sends the data and the control words to the parity-checking block to check the parity.

A common control word format is used in both the incoming and outgoing interfaces. The Figure below describes the fields of the control word. When inserted in the data path, the control word is aligned such that its MSB is sent on the MSB of the incoming or outgoing data lines. A (payload) control word that separate two adjacent burst transfers may contain status information pertaining to the previous transfer, the following transfer, or both transfers. The table below shows a list of control word types. Maximus would ignore reserved control when detected from the incoming SPI-4 interface.

Figures 42 shows the SPI-4 control word format and Figure 43 and 44 contain the detailed description for the SPI-4 control word format.

The Rx state transition table is shown in Figure 45 and is based on the types of control words received. Since the state machine has to decode two 16-bit words in one cycle, the table lists all the possible 16-bit words combinations. The 2-bit ctrl signal is derived from the input DDR signal RCTL. Bit 1 is for the word 1 and bit 0 for word 0. 1 indicates a control word and 0 is a data word.

I_ctrl = Idle control word detected P_ctrl = Payload control word detected

T_ctrl = Training control word detected

Parity Checking

The parity-checking block calculates 4-bit Diagonal Interleaved Parity (DIP-4) and compares the calculated DIP-4 with the one received in the control word. Then it reports the error. The extent of the parity calculation is shown in Figure 46. A functional description of calculating the DIP-4 code is given in Figure 47.

1 Sync FIFO

2 This FIFO is an 8x64 dual-port SRAM used to synchronize the 200MHz
3 data/controls from the incoming SPI-4 to the internal clock. The write operation is
4 controlled by the RX_SM and the read operation is controlled by the per_port control
5 block. The per_port control reads from the sync FIFO as long as it is not empty.

7 SPI TBL

8 The SPI TBL is a table of 64 entries, one for each port. The table contains control
9 information for the incoming data burst for each port.

11 FBL

12 The FBL (Free Buffer List) block is an 84x10 single-port SRAM containing the
13 pointers to the data FIFO that are free to accept data.

15 Data FIFO

16 The data FIFO is a 672x64 dual-port SRAM consists of 84 blocks of 64 bytes. It
17 is used to receive data for each incoming port and only send the data out if it has
18 accumulated 64 bytes or reached end of packet.

20 Calendar

21 The calendar block contains a port calendar conFigured at the startup time. An
22 example of a port calendar is shown in Figure 48.

23 The calendar block receives per-port-basis FIFO status from the segmentation and
24 encodes the status of each port in a 2-bit data structure. Then it sends the status back to
25 the transmit device in a round robin fashion based on the port calendar. The "11" pattern
26 is reserved for in-band framing. A DIP-2 odd parity is sent at the end of each complete
27 sequence, immediately before the "11" framing pattern. See Figure 49. A description of
28 calculating DIP-2 is shown in Figure 50.

30 Memory Description**32 Summary of Internal Memories:**

33 See Figure 51.

34 Interface**35 Interface Ports**

36 See Figure 52.

38 Register Description

39 See Figure 53.

41 Principles of Operation

42 The incoming SPI block receives 16-bit DDR data, 1-bit DDR control and a 400
43 MHz clock from the LVDS pins of the Maximus chip. It first converts the 16-bit DDR
44 data words into 32-bit SDR data words and a 1-bit DDR control into a 2-bit SDR control.
45 The 2-bit control indicates which 16-bit word within a 32-bit word is the control word or
46 data word. Based on the 2-bit control, the Rx state machine starts to extract the in-band

port address, start/end-of-packet indication and error-control code from the 32-bit data path and packs the 16-bit data words into 64-bit data words.

Subsequently, the incoming SPI block outputs the 64-bit data words along with the control information at 200 MHz to the lookup engine. The incoming SPI block also calculates 4-bit Diagonal Interleaved Parity (DIP-4) and reports an error if the calculated DIP-4 is not equal to the one received with the control word.

The incoming SPI block also sends 2-bit FIFO status information back to the transmitting device. This FIFO status information is originated from the segmentation block indicating FIFO congestion conditions for all the ports. The SPI block sends the FIFO status to the transmitting device at the rate of one port per cycle in a round robin fashion based on the port calendar sequence.

Outgoing SPI-4 Interface

The Outgoing System Packet Interface Level 4 (SPI-4) interface transmits 16-bit DDR ATM or Packet over SONET/SDH (POS) data to the receive device (PHY/Switch, depending on whether employed in the ingress or egress application) and receives out of band FIFO status information from the receive device. The outgoing SPI interface gets data and control information from the re-assembly engine, and passes the FIFO status information to the database. Furthermore, the outgoing SPI-4 interface also interfaces with the CPU through a software interface. The CPU can program the registers inside the blocks that comprise the Outgoing SPI-4 Interface.

Features

- Supports SPI-4 phase II spec at 10 Gbps rate.
- Unpacks the 64-bit data at 200 MHz to 16 bit DDR data at 400 MHz.
- Muxes the data and the control onto the SPI-4 bus.
- Diagonal Interleaved parity generation.
- Programmable data transfer burst sizes of multiple 16 bytes with minimum of 80 bytes.
- Passes the FIFO status info to the database to stop transferring.
- Schedules training sequence for data path de-skew.
- Supports maximum 64 ports.

Functional Description

The Outgoing SPI-4 Interface receives 64-bit data and control information from the re-assembly engine at 200 MHz. The Tx state machine generates 16-bit control words, unpacks the 64-bit data into 16-bit data words, and multiplexes the control and data words onto a 16-bit, 400 MHz DDR bus to an external receive device. The Outgoing SPI-4 Interface also passes the flow control information from the receive device to the database block within the Maximus chip. See Figure 54.

Block Description

Tx State Machine

The Tx state machine stores and monitors the length field from the reassembly block. The length field is decremented on every data transfer from reassembly to SPI.

When the length reaches a threshold, the Tx state machine issues a request command to read the next port data and control information out of reassembly block. The threshold level is determined by a value that would allow back-to-back data transfer from reassembly block into SPI without gap on the data bus.

The control word format is common for both incoming and outgoing interfaces. Refer to the Incoming SPI Interface document for details.

The state machine also schedules to send the training patterns for the data path de-skew at least once every pre-defined DATA_MAX_T cycles. The training pattern is shown in Figure 55 from cycles 1 through $20\alpha + 1$. The α is the number of repetitions of the data training sequence. The DATA_MAX_T and α are configured at the startup.

The state transition table is shown in Figure 56. At each state, two 16-bit words and pos_v and neg_v have to be output. And the pos_v and neg_v will be used to generate a DDR control signal. 1 indicates a control word and 0 a data word. If data_valid is set, there are data to be sent. Otherwise the state machine needs to generate one of the control words.

Data Parity Generation

This block calculates 4-bit Diagonal Interleaved parity (DIP-4) and sends the DIP-4 along the control words. The extent of the parity calculation is shown in Figure 57. A functional description of calculating the DIP-4 code is given in Figure 58.

Unpack Engine

This block receives 64 bit data from Reassembly block and samples data at 200MHz clock rate into the SPI out block. See Figure 59.

The 64-bit data and a data valid signal are received following the control info read. The unpack engine unpacks the 64-bit input data into 16-bit words. In order to send 16-bit DDR data at 400 MHz to the LVDS I/Os, the state machine has to come up with two 16-bit words every 400 MHz clock cycle. See Figure 60.

Configuration Register

This block defines the organization of configuration space based on SPI-4 phase 2 specifications.

Port Calendar

This block stores the predefined calendar information to identify each port is ready to send data via the SPI interface. The port calendar contains three fields to indicate the PORT ID, JUMP and VALID. The port ID represents the port number for transmission of data and receiving of status. The port calendar can store up to 64 ports, in the event that less than 64 ports are used, JUMP indicates the end of calendar. VALID indicates if the port is valid or not. If a given port is not valid, the next PORT ID should be used, and the PORT ID marked as not valid should be skipped.

When CPU updates the port calendar after initialization, CPU needs to meet the following conditions:

- CPU can only modify one port at a time.
- When CPU write a port, CPU needs to specify a port entry in the port calendar between port 0 to 63, and the specific port will be modified.

- 1 - When CPU read a port, CPU needs to specify a port entry in the port calendar
- 2 between port 0 to 63, and the specific port will be read out.
- 3 - CPU needs to stop the traffic for a given prior to modifying the port attribute
- 4 (Jump, Valid).

5 See Figure 61.

6 **Port Calendar State Machine**

7 This block waits for rdy_4_nx_port_id from TX state machine to indicate the
8 transmit state machine is ready to accept another port request. Upon receiving
9 rdy_4_nx_port_id, calendar state machine issues the next port ID from the port calendar
10 to the TX state machine. However, before issuing the next port ID, the status bit is check
11 to ensure the SPI FIFO status is in the “Hungry” state for the given port. There is a 64-bit
12 register to retain the value status on a per port basis. The block matches incoming status
13 arriving into SPI against attributes from Port Calendar to identify the given port traffic
14 status.

15 After Port Calendar issues a valid port to the TX state machine, the Port Calendar
16 State Machine pre-fetches the next Port ID from Port Calendar by performing a read. If
17 the TX State Machine issues a read before a valid Port ID is inside Port Calendar State
18 Machine, no_valid_port_id_available is returned to the TX State Machine.

19 **Status Parity Check**

20 This block checks 2-bit Diagonal Interleaved parity (DIP-2) and verifies DIP-2
21 parity bits for signal integrity. Assuming the TSCLK is running at 400 MHz, this block
22 operates at the 400 MHz clock domain and samples the tstat [1:0] twice every clocks.

23 There are two 65-bit registers to hold the FIFO status. The first 65-bit registers
24 hold the status value before parity check (pre_chk_st) and the second register holds the
25 status after parity check (post_chk_st). Once the parity check passes, this block is
26 responsible to copy the status from the pre_chk_st to the post_chk_st register. If there is a
27 parity error, the error is reported into the control status register (spio_err [0]) by setting
28 1'b1 into the bit. This spio_err register is self-cleared upon CPU read.

29 This block sends SPI status from post_chk_st to the Database block. The SPI port
30 status will be sent in spio_dbs_stt [7:0], and status for port 0-7 should be sent when the
31 spio_dbs_sync is asserted, and status for port 8-15 should be sent in the following clock,
32 status 16-23 on the next clock, etc. SPI status spio_dbs_stt [8] should contain the status
33 for the CPU port. spio_dbs_sync should be derived from the global_sync signal.

34 When tstat [1:0] is 2'b01, this indicates “Hungry” and transfer of up to
35 ((MaxBurst2)*(16-byte)) blocks or the remainder of what was previously granted
36 (whichever is greater) may be sent. Since Maximus internal cell size is 64-byte plus
37 header (up to 16 bytes), after SPI-out provides an indication to Maximus to stop
38 scheduling additional cell to memory, SPI-out needs to allow the remainder of current
39 cell to go out of SPI. In other words, reassembly block does not support shipment of
40 partial cell to SPI, each time SPI request for a cell, SPI must guarantee to be able to
41 accept one complete cell. Therefore, the MaxBurst2 value should be programmed to
42 minimum of 5. Once programmed to 5, when the “Hungry” indication is received, the
43 remaining bytes in a given cell can still be shipped out. To notify Maximus to stop
44 scheduling the next cell of data is the spio_dbs_stt bit from SPI-out to Database block.

Consequently, spio_dbs_stt should be set to one when “Hungry” is received for a given port.

This interface between SPI-out and database drives the status of the 64 ports in the output SPI and the CPU port. It is a continuous interface. Every 8 cycles of 200MHz the status indication of the 65 ports is updated to a 65 bits register. The first cycle of the eight is marked with a specific signal (sync signal).

Figure 62 shows the interface:

A synchronization FIFO is used between the 400 MHz and 200 MHz clock domain.

This block detects the synchronization pattern from tstat [1:0] by comparing the incoming tstat[1:0] against the framing pattern. Also, the port calendar state machine checks training sequence to see if there is any skewing on FIFO status and control lines. To check the training sequence for skew, this block needs to implement a counter for internal tstat, called sync [1:0], which cycles the nominal value of tstat [1:0] with respect to Maximus internal clock. The value of sync is compared against the value of tstate [1:0] (external) from the SPI pad and the difference between the two pairs of signals determines if there is any skew of control signals. See Figure 63.

TX FIFO

This block implements a 4x16-bit dual port SRAM to synchronize Maximus and SPI clock domains. On the incoming side, the write port operates at 200 MHz and pushes data into the TXFIFO when data_push is asserted by TX state machine. On the outgoing side, the read port operates at 400 MHz DDR and pops data into the tdat [15:0] when data_pop is asserted by TX state machine.

TX Status FIFO

This block synchronizes status from 400 MHz to 200MHz domain using 4x16-bit dual port SRAM.

Memory Description

Data Structure for Port Calendar

See Figure 64.

- Memory type: single-port SRAM
- Memory configuration: 64 x 8 bit
 - o Port ID: hold port identification value programmed by the CPU.
 - o Jump: indicates the end of the port calendar. When Jump is asserted, the given location is the last entry on the port calendar. Upon detecting the assertion of jump, the port calendar should recycle the port ID from the top of the port calendar.
 - o Valid: When valid is asserted, it indicates the given port hold a valid connection. When valid is de-asserted, the port calendar will skip the de-asserted port and go to the next port ID on the port calendar.

Data Structure for TX FIFO

- Memory type: dual-port SRAM

- Memory configuration: 4 x 16 bit
- The TX FIFO is used to synchronize data between different clock frequencies.

Data Structure for Status FIFO

- Memory type: dual-port SRAM
 - Memory configuration: 4 x 16 bit
- The Status FIFO is used to synchronize data between different clock frequencies.

Interface

See Figure 65.

Timing Diagram

Timing information between SPI-out and Reassembly, see Figures 66 and 67.

Register Description

Start-up Parameters, see Figure 68.

Software Interface

The software interface programs Outgoing SPI with a CPU. All Outgoing SPI's control registers and memories can also be written and read for test purposes through the software interface. The following sections list all the control registers and all the Outgoing SPI software interface commands. Please see a dedicated document for a detailed description of the software interface.

Registers

See Figure 69.

Commands

Read Outgoing SPI Start-up Parameters.

See Figure 70.

Write Outgoing SPI Start-up Parameters

See Figure 71.

Read Outgoing SPI Start-up Parameters

See Figure 72.

Write Outgoing SPI Start-up Parameters

See Figure 73.

1 **Read Outgoing SPI Start-up Parameters**

2 See Figure 74.

3 **Write Outgoing SPI Start-up Parameters**

4 See Figure 75.

5 **Write PORT Parameters after Initialization**

6 See Figure 76.

7 **Read PORT Parameters after Initialization**

8 See Figure 77.

9 **Read PORT Full Parameters**

10 See Figure 78.

11 **Write PORT Full Parameters**

12 See Figure 79.

13 **Read Control Status Register**

14 See Figure 80.

15 **Write SPIO_TSTMUX_SEL**

16 See Figure 81.

17 **Read SPIO_TSTMUX_SEL**

18 See Figure 82.

19 **Reserved Op-codes**

20 All op-codes other than listed above are reserved.

21

22 **Principle of Operation**

23 The outgoing SPI block receives FIFO status from the Incoming SPI. The
 24 beginning of the status is found based on the framing sequence. The beginning of the
 25 status from SPI interface is mapped to the port calendar to look up the status of each port.
 26 The status has to be synchronized to the internal 200 MHz clock and passed to the
 27 database block, and internal state machine to select a port that is available for data

transmission. Base on the status information, outgoing SPI block requests data from re-assembly block on a given port. Outgoing SPI unpacks the port data from re-assembly, delineation data between different ports and check for data error using control information from re-assembly. The outgoing SPI block calculates the 4-bit Diagonal Interleaved Parity (DIP-4) and combines those four bits with the control information to form a 16-bit control word. At every clock cycle of the 400 MHz clock, the outgoing SPI unpacks the 64-bit data into 16-bit data words and inserts control words by multiplexing onto 16-bit DDR bus at 400 MHz. If the re-assembly's control FIFO is empty, the outgoing SPI block sends idle control words onto SPI data bus. During regular interval, the outgoing SPI also sends training control word to provide skew detection.

Lookup Engine

This block receives ATM/Packet data and control from the SPI-4 receive interface block. It looks up a Flow ID and a set of parameters from a hash table in external SSRAM and forwards the data flow with the parameters to the segmentation engine.

Features

- Supports up to 1M flows.
- Supports up to 64 input ports.
- Operates at 200 MHz/5ns cycle time.
- Tables and configuration registers loaded via CPU Interface.
- Label extraction is programmable on a per port basis.
- Accepts packet SOP cells of 64 Bytes only.
- Accepts ATM cells of 52 Bytes only.
- Direct header stripping is programmable on a per port basis.
- Hash generator creates 4 hash values or slots in parallel from a key using 4 separate CRC polynomials as hash functions.
- Two bins per bucket or record to minimize collisions.
- There are two record memories for up to 8 read operations per a tag lookup.
- Hash table lookup constrained to resolve within 8 read operations.
- Overflow CAM of 32 X 72 entries for lookups that fail to resolve in 8 read operations.
- Fully pipelined lookup operations.
- Supports a maximum of two 8 X 72 M SSRAM. The normal size is two 2 X 72 M SSRAM.
- Supports the following flows with the following headers:
 - MPLS with PPP encapsulation.
 - LC-ATM that is VPI/VCI.
 - Switch headers.
 - FIDS.
- Operates in ingress and egress modes.

Functional Description

See Figure 83.

1 The lookup engine does label extraction, label hashing and table lookup for each
2 data transfer that is SOP with a header. The type of traffic is unique to the incoming port.
3 There will be no more than one type traffic on a port. The default type will get a special
4 FID. The default FID is stored in a memory and it is specific to the input port. The lookup
5 engine can strip the switch header of incoming cells.

6 **Block Descriptions**

7 See Figure 84.

10 **CPU Hash Generation**

11 The CPU through the interface can request a hash generation and get a hash index
12 as a response. It will be used to set or tear down connections.

14 **Port Configuration**

15 The memory has two uses. It specifies the type of traffic on that port. Also, it has
16 the Traffic parameters. The CPU either writes the parameters or the logic after lookup. If
17 the CPU writes the parameters, then they are not changed until the traffic type is. They
18 will be used in default traffic. If the traffic is direct, then the rest of the parameters in the
19 memory for that port are not used. If the traffic is ENCAP type, then another set of
20 registers are used. When the port traffic is none of the above, then the parameters values
21 such as the type and FID are looked up from the record memory for SOP traffic. Then at
22 the last stage they are driven on the interface bus. At the same time they are written
23 immediately in the port configuration memory. For none SOP cells, the parameters are
24 read from the port memory at the last stage and then driven on the bus. This will avoid
25 data coherency problem in case the cells for the same port are contiguous. For ATM cells
26 the SOP and EOP cell is the same.

28 **DATA Coherency**

29 The parameters on none SOP cells are read near the last stage of the data in order
30 to avoid driving stale values on the interface bus. The CPU accesses to modify the
31 memory are not tracked. The parameters used are the ones that are lookup or stored in the
32 port memory.

33 If the SOP cell has a match on the last lookup, the parameters are written in the
34 port memory near the end of the lookup cycle. If the next data cell is for the same port,
35 then the parameters are not written yet in the corresponding memory location. Therefore,
36 if the cell is not SOP, the port memory is read near the end of the pipeline of the cell
37 lookup. See Figure 85.

39 **CPU Accesses**

40 The CPU accesses are allowed between lookup sequences when the memory bus
41 is not busy.

43 **External SSRAM**

44 The Flow ID tables are implemented in off chip rams using 200MHZ, ZBT
45 SSRAMS. These rams provide single cycle access with no 'dead cycle' penalty when
46 mixing reads and writes. The table will contain a maximum of 8M entries. There is a

programmable register that will determine the size of the memory. The hash result will be trimmed to a number of bits that is programmable. There are two SSRAM channels provided. The usual size is 2M entries per channel. The memory can be accessed directly without hash for a read and write operations.

SRAM INITIALIZATION

The internal as well as the external ram will be initialized to

HEADER Selection

The header depending on the type has different sources. It can be transmitted with the traffic as in the case of direct or the special type. It can be found in the port memory as in the case of the default. Last, it can be found in record memory. See Figure 86.

Hash FIFO and Record Reading

See Figure 87. The data is stored until the port table is accessed. The label is generated. The data will continue to be clocked in a FIFO until the key is hashed and the index is calculated.

Memories

Input Port Parameters Memory

Each port will have a parameter entry in that table. The "type" parameter will determine if the rest of the entries are used. If the type is "000", then the rest of parameters of that memory entry are not used except for the header strip and tag extraction. The "HD STRP" indicates the number of bytes to strip from the header. The "INT HDR EXTRACT" indicates the number of bytes that represent the header. The header is at the end of the header stripped that is stripped. If the type is default, then the header to strip is zero. If the header to strip is not equal to zero, then the header will be stripped as in the case of the direct. The same argument is for the tag extract. The header will contain the FID, CLP, EFCI, and the OAM bits. The bits will be extracted and used in lieu of the look up data. The number of bytes of the header that are stripped and the parameter location is specified in the table. If the type is "001", then it is a default type and the rest of the parameters of that memory entry are used. If the type is "001", then the incoming traffic will be treated differently. The traffic is special and it is treated later in the specifications in detail. Again there are no hashing or lookup operations. If the type is different than the above three, then the tag is extracted, hashed and the parameters are looked up from memory. There are 64 entries in the table because the device can support a maximum of 64 ports. When there is a lookup operation, the parameters that are looked up are stored in the port memory and are used on the traffic until the next SOP cell. Infields where the hardware and the software perform the write operations, the software does it for the default traffic only, the other types of traffic are updated by the hardware. See Figure 88.

The CPU writes the data in the port table entries for the default or the direct type. The logic will write the parameters for the port if it the type is not any of the above. The SOP cell will trigger the lookup. The rest of the cells on that port will use the table until the EOP or the next SOP. Then the parameters will change again.

1 In L2 encapsulation, if the CRC bit is set, then assert the CRC signal to the segmentation
2 engine. See Figures 89 and 90.

4 **Table Lookup Memory**

5 The index or hashed key provided by the hash generator is used as an address.
6 The number of bits that are used from the hashed key or index is programmable and
7 depends on the size of the SRAM. The address is a base one and it is used to access four
8 entries in the hash table located in external SSRAM. When these locations are read, the
9 hash key is compared to the content of the memory. If the first location does not match,
10 then there is collision and the second one is compared the same way. The process is
11 repeated for all four entries. If there is no match and collisions resulted from the compare
12 operations, then the process is repeated for the second hash result or index and so forth. If
13 at the end there are only collisions, then the CAM is accessed and the process of
14 comparison is repeated. If there is a hit, then in memory location there is Flow ID (FID),
15 the QOS, and the rest of the control bits. Each hash table entry requires 72 bits. The bit
16 definitions are summarized in Figure 91. The output port number is used to access a
17 register. If the register bit corresponding to that port is set, then the EOP cell will have
18 the error or discard flag set.

20 **Hash Data Memory**

21 See Figure 92. The data FIFO is a 15 X 104 bit wide FIFO. The FIFO
22 compensates for the pipeline delay of the lookup path. It also provides buffering that is
23 needed when the segmentation engine is not ready to accept any data. The data receiver
24 block writes the FIFO. The data transmitter block reads the FIFO. The segmentation
25 engine will assert a ready signal when it can accept data. When the ready is asserted, then
26 the data is pushed to the segmentation engine.

28 **Interface**

29 See Figures 93 and 94. When the segmentation engine asserts the
30 "ALMOST_FULL_SIGNAL", the lookup engine will finish transferring the cell of eight
31 transfers. Then at the cell boundary, the lookup engine will terminate transmission. At the
32 same time the lookup engine will signal the SPI interface to terminate pushing data to the
33 look up engine. The pipe of tag lookup operations is not violated. There are no lookup
34 operations when the traffic flow is suspended.

36 **CPU Interface**

37 There are signals that are generated by the CPU interface and are used to control
38 the enable and disable signals. See Figures 95 and 96.

40 **Software Interface**

41 See Figures 97 and 98.

42 The default value for the memory size is "001" which corresponds to 2 MEG.

44 **FID Memories Select**

45 See Figures 99 and 100. The default value for chip select is "00".

Register Description

The CPU interface doesn't support burst capability, instead, any read or write command will terminate with one address location; data bus is still 32 bits. Lookup Engine block has a fixed size of addresses allocated to it and all accesses to any internal/external memories or register is done via that specific address space. This method require the CPU Interface manager to have a mapping of the address spaces per block so that an incoming command will be redirect to the correct block.

Features

- Data bus of 32 bits.
 - Single address access (not a burst access)
 - All external memories and most of the internal memories are accessed via indirect accesses through registers
 - Each block has dedicated addresses for its internal registers of up to 64 registers.
- A command drives address and command (Rd or Wr) to the address FIFO and data if it's a write command. Each block starts it's own CPU access by POPING the address/command from the FIFO. It means that any CPU transaction starts when the address/command FIFO is not empty.

Address space

The address space is 64 entries only. The internal registers directly mapped to these addresses and external/internal memories are indirectly mapped. The address spaces per lookup engine are shown in Figure 101.

Direct/Indirect Accesses

A direct access is used to access. If there are internal memories or registers, which are longer then 32 bits, then read or write to these locations will be executed in more then one access. For example an entry in internal memory of 50 bits width (or a register with the same length) will be accessed in two CPU commands first access bits [31:0] and then [49:32]. Each register or small internal memories have their own address, which is unique to the entire chip set.

An indirect access is used to access internal and external memories. Each block has a set of registers dedicated for the indirect access. One register contains the address and the command (Rd, Wr, Init etc). Few other registers contain the data for a write command. Writing to the COM/ADD register triggers the execution of the command, once the block completed the execution of the command, the COM portion of the command is cleared by the block. When a read command occurs the data is written to the registers and then the CPU reads the registers to get the data.

Registers

See Figure 102.

1 2 **Commands**

3 **NOP**

4 See Figure 103.

5 **Read LUT Memory1 (72 bits)**

6 See Figure 104.

7 **Write LUT Memory1 (72 bits)**

8 See Figure 105.

9 **Read LUT Memory2 (72 bits)**

10 See Figure 106.

11 **Write LUT Memory2 (72 bits)**

12 See Figure 107.

13 **Read Port Memory (44 bits)**

14 See Figure 108.

15 **Write Port Memory (44 bits)**

16 See Figure 109.

17 **Setup connection command for KEY**

18 See Figure 110. This command sets up a connection by accessing external
19 memories and initializes the corresponding locations to the proper parameters. It is a
20 write command.

21 **Tear-Down connection command for KEY**

22 See Figure 111. This command tears down a connection by accessing external
23 memories. It is a write command

1 GET HASH FID

2 See Figure 112. The response will be in the register. See Figure 113. This
3 command tears down a connection by accessing external memories. It is a write
4 command.

5 Read CAM Memory1 (72 bits)

6 See Figure 114.

7 Write CAM Memory1 (72 bits)

8 See Figure 115.

9 Read CAM Memory2 (72 bits)

10 See Figure 116.

11 Write CAM Memory2 (72 bits)

12 See Figure 117.

13 INIT ALL MEMORIES

14 See Figure 118.

15 INIT MEMORY 1 & CAM 1

16 See Figure 119.

17 INIT MEMORY 2 & CAM 2

18 See Figure 120.

19

20 General commands

21 The following commands are combinations of the above simple commands. A
22 user command like this should be broken to several Read or Write commands in the
23 driver.

24

25 Initialized External memories

- 26 • Write "0" in all Record Memory 1.
27 • Write "0" in all Record Memory 2.

1

2 **Initialized Internal memories**

- 3
- Write "0" in all Port memory.

4

5 **Principle of Operation**

6 The lookup engine has to perform a port lookup before the FID lookup operation is
7 performed. The memory is initialized and the valid bit is set for that FID. IF the valid bit
8 is not set, then there will be no match. The lookup engine has the port of incoming traffic.
9 As a result the following is done:

- 10 1. Read the port parameter entry by accessing the port parameter table.
11 2. Get the LOOKUP TYPE.
12 3. Extract the tag.
13 4. Perform the stripping, hashing, or lookup operations according to the type.

14 Once the parameter entry is read, then the lookup type is known. If the lookup type is
15 *direct*, then the lookup engine will extract the FID, QOS, CLP, EFCI, and OAM if any
16 and push all the control and data to the segmentation engine. If the traffic is *default*, then
17 the parameters from the port table are used. If the traffic is *L2*, then the operation is
18 handled as described in another section. If the traffic type is none of the previous ones,
19 then the tag is extracted according to the type of traffic and then it is hashed. The result of
20 hashing will be used to access the memory to get the rest of the parameters if there is a
21 match. The memory content of the tag is compared to the hash tag. If they are equivalent,
22 then the rest of the memory content is used for that FID, QOS, CLP, EFCI and OAM. If
23 the tags do not compare then there is collision and the next hash result is used. There are
24 four hash generators that will generate 4 indices. Two indices will be used for one
25 memory and the other two for the second memory. Each record holds two bins. The total
26 is 8 bins that might collide. If there is no match but collisions, then the overflow CAM is
27 accessed. If there is collision or no record, then the packet or the cell is marked for CPU
28 and a register is used for the FID.

29

30 **CPU FIDS**

31 If the keys stored in the memory locations do not match the hashed key, then the
32 CAM content will be compared and the if it fails to compare, then the traffic will be sent
33 to the CPU. There is a register that holds the FID and the rest of the bits necessary for the
34 data traffic.

35

36 **TAG Extraction**

37 The lookup engine requires the minimum burst size of 16 bytes for the key
38 extraction. The hardware extracts a contiguous number of bytes starting at an offset
39 according to the traffic type of that port. The type of label to be extracted is determined
40 by the port number. The location of the label is dependent of the type of traffic and is
41 hard coded in the hardware.

42

1 **ATM TYPES**

2 If the type of the traffic is ATM, then the CLP, EFCI, and OAM bits are extracted
3 and driven to the rest of the device. See Figure 121. Extract the above bits from header
4 and drive to segmentation.

5

6 **TAGS**

7 The following Figures show the data formats for each of the supported protocols.
8 The specifications are not hard coded. The port table will be used for the header's
9 location. Each field is 8 bits wide with 12 bits VPI and 16 bits VCI field, for a total of 28
10 bits lookup tag. This 28 bits lookup tag is located in the 1st 64-bit word, bits 32 thru 63.
11 See Figures 122 and 123. The PPP protocol is part of the link layer header. The tag is bit
12 60:41 of the first 64 bits word.

13

14 **Port Types**

15 The first lookup is to read the port parameters. If the type is "000", then the type
16 is direct flow. There are no hash or lookup operations. The FID will be extracted from the
17 incoming traffic itself. If the type is "001", then it is a default traffic. The default
18 parameters are stored in the table itself. There are no hashing operations or lookup in this
19 case. The third special case is the L2 type. The incoming header will be used to find the
20 encapsulation header. There will be no hash and lookup operations either. The rest of the
21 lookup types require a hash and lookup operation. The port table will specify the type of
22 lookup that the lookup engine needs to perform on the incoming data. The lookup
23 operation will assign an FID to the data stream or packet or flow, except in the bypass
24 type. The FID is used only internally. See Figure 124.

25

26 **Direct Header Format**

27 See Figures 125-7. The SOP and EOP have to be extracted from the header. The
28 SPI interface will also assert the SOP and EOP signals. The lookup engine will determine
29 the set that uses and asserts to the segmentation engine from the type of the packet. In
30 egress mode, the device uses the header set of SOP and EOP signals depending on the
31 type. For example, if the switch fabric is a packet switch fabric, then the lookup engine
32 will use the SPI interface set. If the switch fabric is a cell one, then the lookup engine will
33 use the set parsed out of the header. When the header is stripped, then the number of
34 valid bytes has to be adjusted accordingly. The data has to be packed in 64 bit transfers to
35 the segmentation engine.

36

37 **Special Header Addition**

38 **If the type is "010", then the flow will be handled differently. The incoming frame**
39 **structure will look like the following:**

40 See Figure 128. The CID and the CLASS will become the new FID of the flow.
41 The rest of the incoming header will be used as the internal header. If the "SPCL LK" is

set, then the encapsulation header is added. If it is not set, then the header is extracted and no special header is added. The memory will have the structure depicted in Figure 129.

The memory location will be read and the length will indicate the number of valid bytes that have to be appended before the data as a special header. The header can be up to 16 Bytes long. The rest of the bits such as the EFCI, CLP and the OAM are driven, too as out of band signals. See Figure 130. The byte valid from the SPI interface to the segmentation engine has to be adjusted for the fact that the extra bytes have been added. The data has to be packed if the number of byte is not on 64 bit boundary.

Hash Index Generation

The extracted label is concatenated with the input port to produce the key. Keys may be padded out to 32 bits if they are less than 32. The key is hashed in hardware using the irreducible ETHERNET CRC32 polynomial listed below and three irreducible trinomial polynomials:

1. $G(X) = X^{32} + X^{26} + X^{23} + X^{22} + X^{16} + X^{12} + X^{11} + X^{10} + X^8 + X^7 + X^5 + X^4 + X^2 + X^1 + 1;$
2. $G(X) = X^{31} + X^7 + 1;$
3. $G(X) = X^{31} + X^6 + 1;$
4. $G(X) = X^{31} + X^3 + 1;$

The key is hashed 32 bits at a time, requiring 1 cycle for all keys. See Figure 131. The type will determine if there is a need for hashing. The tag will be extracted in the first clock because it is available in the first 64 bits transfer. The key will be formed. The key will then be hashed using four different hashing functions. The hashing is done in parallel. The first hashed key or index will be used as a base register to read two memory locations. As soon as the first location is read, then the key that is stored in that memory location is compared to the hashing key. If there is a match, then the read operations are halted. The content of the memory of the matching key will be used. If there is no match, then the next location content will be compared and so forth. The second index will be used to read the third and forth records if the first record of two bins fails to get a match. Then the third and the forth indices are used for the second memory. If all four indices fail to get a match, then the CAM is accessed. If the CAM contents do not match, then the traffic is forwarded to the CPU.

Timing Diagrams

See Figure 132. The transfer size is assumed to be 64 bytes for up to eight transfers for SOP SOB data traffic.

Segmentation

The segmentation engine segments the incoming packets into 52-byte ATM, 48 bytes AAL-5, or switch cells and stores them in memory as cells of 64 bytes. The segmentation engine adds the necessary number of pad bytes to complete the cell.

Main Features

- Segments incoming data from SPI-4 on per port basis, supports maximum of 64 ports
- Segments incoming data from CPU interface

- 1 - 64K bytes maximum packet
- 2 - Segments, maps or de-maps incoming data into either ATM/AAL-5 cells or
- 3 switch cells
- 4 - Stores segmented cells in memory as cells of 64 bytes
- 5 - Supports the traffic depicted in Figure 133.
- 6 - Types of segmentation and applications for various traffic types.
- 7 - Backpressures the scheduler in the egress side from the per flow Q in the
- 8 ingress side via status cells through a serial interface to segmentation engine.
- 9 - Calculates CRC-32 for packets segmented into AAL-5 cells.
- 10 - Calculates new CRC-32 for L2 header.
- 11 - Supports programmable maximum transfer unit up to 16K bytes for packet
- 12 data.
- 13 - Each port has a 32-bit counter to count packets coming from that port.
- 14 - Each port has a 48-bit counter for data through that port.
- 15 - Tags bad packets that they can be discarded later.
- 16 - Internal memories to store parameters and data.
- 17 o Dual-port SRAM
 - 18 ▪ Control - 65 x 39 bits to store packet length, and cell byte count
 - 19 ▪ Data - 640 x 64 bits to store data
 - 20 ▪ CRC - 65 x 32 bits to store the partial CRC for each port
 - 21 ▪ Statistic - 65 x 80 bits to store packets and data received for 64
 - 22 ports, and one for CPU
- 23 o Single-port SRAM
 - 24 ▪ Buffer - 80 x 7 bits to maintain the free buffer list for data
 - 25 SRAM
 - 26 ▪ Queue - 80 x 70 bits to queue information for the complete
 - 27 cells
 - 28 - Operates at 200 MHz

Functional Description

See Figure 134. There are three input data paths to the segmentation engine: data from SPI-4, CPU, or serial interface. There are two output data paths from the segmentation engine: to memory manager, and serial interface to the scheduler.

See Figure 135. The segmentation engine receives data from the incoming SPI-4 during normal operation. Based on the type of application, the segmentation engine segments, de-maps, or maps the incoming data to either ATM/AAL-5 or switch cells. As shown in Figure 133 there could be five different types of segmentation (as opposed to types of application) performed by the segmentation engine. The following explains the operation of the segmentation engine for each of these types.

TYPE 1 Segmentation

See Figure 136. Type 1 segmentation is for application type 0, 4, 8, 9, 10, 11, 12, and 14. The segmentation engine receives the 64-byte cells from the look up engine and queues the cells to the output control block when the entire cell has been received.

TYPE 2 Segmentation

See Figure 137. Type 2 segmentation is for application type 1 and 5. The segmentation engine receives the 64-byte cells from the look up engine and queues the cells to the output control block when the entire cell has been received. The cells come in as 52 bytes ATM data with 12 bytes of pad and output as 48 bytes of AAL5 data with 16 bytes of pad to the memory manager block. The segmentation engine removes the 4 bytes ATM header and appends 4 bytes of pad to the end of the cell.

TYPE 3 Segmentation

See Figure 138. Type 3 segmentation is for application type 2 and 13. The segmentation engine receives the 64-byte cells from the look up engine and queues the cells to the output control block when the entire cell has been received. The cells come in as 64 bytes packet data and output as 48 bytes of AAL5 data with 16 bytes of pad to the memory manager block. The segmentation engine segments the incoming data into 48 bytes, appends 16 bytes of pad, and passes the control information for the completed cell to the output control block. CRC-32 is calculated when the cell is being sent out to the memory manager block.

TYPE 4 Segmentation

See Figure 139. Type 4 segmentation is for application type 3 and 6. The segmentation engine receives the 64-byte cells from the look up engine and queues the cells to the output control block when the entire cell has been received. The cells come in as 64 bytes packet data and output as 64 bytes of AAL5 like data to the memory manager block. The segmentation engine segments the incoming data into 64 bytes and passes the control information for the completed cell to the output control block. The segmentation engine calculates CRC-32 and appends an AAL5 like trailer at the end. However, CRC-32 is calculated when the cell is being sent out to the memory manager block.

TYPE 5 Segmentation

The segmentation engine actually is not segmenting data in type 5. When the ingress segmentation engine receives a status cell from the egress PFQ through the serial interface, it queues that status cell to the memory manager like any other cells. The status cell is sent through the switch to the egress segmentation engine, and the egress segmentation engine sends the status cell to the ingress scheduler through a serial interface. See Figure 140. The serial input and output data paths provide a way for the per flow Q on the egress side to back pressure the scheduler on the ingress side from sending data.

CPU injected packets

The CPU can inject data through software command as described later in section 1.6.2. When the segmentation engine detects the CPU inject packet command it sends the cell out in the next control pop from the memory manager. However, there can only be no more than one CPU cell sent within 64 cells time. The segmentation engine starts a cell counter whenever it sends a CPU cell, and no more CPU cells will be queued out until the counter reaches 64.

As shown in Figure 141, the CPU can write and read the internal registers inside the segmentation engine through the CPU interface. The CPU can also input packets into the segmentation engine through the CPU interface.

Block Description

See Figure 142.

Description of INPUT CONTROL block

In normal operation, when the segmentation engine has room in its FIFO for incoming cells it asserts a FIFO not full signal to the look up engine and waits for a valid signal from the look up engine to indicate there is valid control information on the bus. When there are valid controls, the input control block performs read from various internal memories:

- a) Reads from the FBL SRAM to get two pointers to the data SRAM,
- b) Uses the port address to read the SEG_TBL SRAM to get the per port parameters.
- c) Uses the port address to read the STATISTIC SRAM to get the statistics for that port,

The per port parameters (SOP, EOP, etc.) from the bus are compared with those from the SEG_TBL, and error conditions are set accordingly. The port address from the bus is compared to the CUR_PID, which stores the last port address processed, to see if the incoming burst is from the same port address. If so, it means the data burst is the continuation to the one previously processed, and the control information in the internal registers are used instead of those read from the memories. Otherwise, it means the incoming data burst is for a new port address and new controls are to be used, and the controls for the last processed one are to be stored before the new ones are loaded. Detailed information is discussed later in the timing flow.

The input control block has two sets of control parameters, one for the last processed port address, i.e. the current parameters, and one for the port address from the look up engine that comes with the valid signal.

The controls, statistics, and the two pointers are then passed to the data block with a start signal to start the data state machine.

The input control block also receives FIFO status from the FBL and QUEUE blocks and generates the FIFO not empty signal to the look up engine.

Description of FBL block

The FBL (Free Buffer List) block keeps track of the open entries in the data SRAM. After power on reset, the CPU needs to write the init bit in the segmentation register to start the self-initialization of the free buffer list. The self-initialization takes a little over 80 cycles, and no data is allowed through the segmentation engine block until the initialization is done.

The FBL block receives the read strobe from the control block. For each read strobe the FBL block reads two locations from the FBL SRAM and passes the two pointers back to the control block.

The FBL block receives the write strobe and the return pointer from the data block. It writes the return pointer to the write address when there is no read pending.

The FBL block generates a full status to the control block if there are no more free entries in the data SRAM.

1 **Description of FBL SRAM block**

2 This is an 80 x 7 bits single-port internal SRAM to store pointers to the free
3 memory blocks in the data SRAM. The CPU needs to set the init bit to start the self-
4 initialization after reset.

6 **Description of SEG_TBL SRAM block**

7 This is a 65 x 19 bits single-port internal SRAM to store the parameters for the 64
8 ports plus CPU. This RAM is accessed twice by the control block for each cell comes in
9 for a particular port, one read to retrieve port parameters and one write to return
10 parameters for that port.

12 **Description of STATISTIC SRAM block**

13 This is a 65 x 80 bits dual-port internal SRAM to store number of packets (32
14 bits) and data bytes (48 bits) received on the 64 ports and CPU interface for statistic
15 purpose.

17 **Description of DATA PIPELINE block**

18 The DATA PIPELINE block delays the incoming data and byte valid to match the
19 delay for control signals through the INPUT CONTROL block.

21 **Description of DATA block**

22 The data block receives the control information from the input control block and
23 decodes the types of application to be performed by its input state machine. The input
24 state machine starts when the data block receives the start signal from the control block.
25 The input state machine has five different branches based on the segmentation types
26 decoded from the application types. For each of the branch the input state machine
27 performs the following operations and writes the data into the data SRAM:

- 28 a) Add 8 bytes of pad
- 29 b) Add 16 bytes of pad
- 30 c) Segment incoming data into 48 bytes and add 16 bytes of pad
- 31 d) Segment incoming data into 64 bytes
- 32 e) Pass through the 64 bytes

33 The main function of this input state machine is to accumulate a 64 bytes cell for each
34 port address. When there is a completed cell, the input state machine generates a control
35 signal to queue the control information for that particular cell into the QUE FIFO in the
36 queue block. At the same time, the data block writes the statistics back into the STAT
37 SRAM and sends the parameters back to the input control block.

38 There is a data byte counter which increments for each incoming data bytes to count
39 for the 64 bytes cell. For a cell with EOP in it, the counter value is sent to the memory
40 manager indicating the number of bytes valid in the EOP cell.

41 For AAL5 or like data cell, the data block also check to see if a new cell is needed for
42 the trailer. If needed so, the data block writes the data/trailer to the buffer entry indexed
43 by the second pointer. Otherwise, the data block returns the second pointer to the end of
44 the free buffer list.

45 There is also an output state machine, which starts when the queue block requests to
46 de-queue from the data block. This output state machine receives the pointer to the data

1 SRAM and the control information to be processed on the cell. The output state machine
2 reads from the data SRAM with the buffer pointer and output the data to the memory
3 manager. When the data block finish sending data to the memory manager it returns the
4 pointer back to the free buffer list.

5 If the control information indicates a need to generate CRC-32, the data block reads
6 the CRC SRAM with the port address from the control information and load the partial
7 CRC into the CRC-32 engine. At the end of the cell and not EOP, the data block writes
8 back the partial CRC into the CRC SRAM. If the control indicates a need to calculate
9 new CRC for the L2 header, the output state machine replaces the existing CRC bytes
10 from the packet with the ones calculated with the new L2 header, and the new values are
11 sent through the CRC engine to calculate for the AAL5 trailer.

12 There are two statistic counters in the data block, one for the number of packet
13 received and the other for the number of data received. These statistics are kept on a per
14 port basis, where the data block load the information from the STAT SRAM with the
15 load signal from the input control block.

16 17 **Description of DATA SRAM block**

18 This is a 640 x 64 bits dual-port internal SRAM to store the incoming data as the
19 segmentation engine completes segmenting them into 64-byte cells. The memory is
20 allocated into 80 blocks with each block has eight entries, total of 64 bytes per block.

21 22 **Description of CRC_GEN block**

23 This block generates the CRC that covers the PDU (Protocol Data Unit) in the
24 AAL-5 format. The polynomial used is CRC-32 as follows.

$$G(x) = x^{32} + x^{26} + x^{23} + x^{22} + x^{16} + x^{12} + x^{11} + x^{10} + x^8 + x^7 + x^5 + x^4 + x^2 + x + 1$$

25
26 Since the internal data bus has 64 bits, the CRC-32 is calculated in two portions:

- 27 a) CRC-32 engine with 64-bit data input
- 28 b) CRC-32 engine with 32-bit data input

29 From the beginning of the data, CRC-32 engine with 64-bit input data bus is used
30 all the way to the end of the data, i.e. right before the trailer for AAL5. Then it loads the
31 remainder from the CRC-32 engine with 64-bit input data bus to the CRC-32 engine with
32 32-bit input data bus and calculates the CRC with the four MSB bytes of AAL5 trailer.
33 The CRC_GEN block reads the partial CRC with port ID as the address to the CRC
34 SRAM. The CRC_GEN block loads the partial CRC into the CRC engine and starts
35 calculating as the data are going out to the memory manager. When the cell is sent to the
36 memory manager block, the CRC_GEN block writes back the partial CRC into the CRC
37 SRAM, again using port ID as the address. At the end of a packet, the CRC is muxed out
38 as part of the AAL-5 trailer. For L2 header, CRC is calculated for the packet and is used
39 to replace the CRC bytes in the packet before calculating the CRC for the AAL5 trailer.

40 41 42 **Description of CRC SRAM block**

43 This is a 65 x 32 bits dual-port internal SRAM to store partial CRC for each port.
44 The data block reads and loads the partial CRC if the traffic type requires CRC-32
45 calculation. The data block writes back the partial CRC when the completed cell is sent to
46 the memory manager.

Description of OUTPUT CONTROL block

The OUTPUT CONTROL block sends the available signal to the memory manager to indicate that there are completed cells available to transfer. The memory manager replies with the control pop signal if it is ready to receive the completed cells. When there is control pop from the memory manager, the output control block reads the QUEUE FIFO for the control information for that cell and passes them to the data block. The data block performs read to the DATA memory as described in the data block above. The output control block sends the control information lining up with the output data. The output control has to work in compliance with the global sync, it only asserts the available signal in second clock slot from the global sync (see detailed timing in later section). The memory manager is to responds in the third clock slot, and the output control block generates the pop signal to the QUEUE FIFO in the following clock. The output control block also generates the read signal to the data memory in the same clock slot.

If there is status cell pending, the output control block sends the status cell using one of the regular port accesses from the global sync, and the skipped cell takes over the CPU access during that particular global sync. This means CPU access is skipped for the global sync that has a status cell to send. The output control block returns to its normal port access in the following global sync.

Description of Q_FIFO SRAM block

This is an 80 x 70 bits single-port internal SRAM to store control information for cells queued to the memory manager block.

Description of CPU_INTF block

The CPU_INTF block is to interface with the external CPU. It supports CPU direct access to the internal registers and indirect access to the internal memories. CPU can also injected packet data through this interface. The CPU can inject packet/data through use of CPU inject data command, however, it can only injects a 64 bytes cell one at a time. As soon as there is a complete cell, the CPU_INTF block asserts a signal to the output control block to output the cell in the CPU access.

Description of PFQ_SEG INTF block

The PFQ_SEG INTF block receives data from the serial interface between per flow Q and the segmentation engine. The PFQ_SEG INTF block shifts the serial data into a 64-bit shift register and load the output of the shift register into an internal register. The control information is extracted from the serial data. When the PFQ_SEG INTF block receives the entire status cell, it queues the completed cell into the queue FIFO.

Description of SEG_SCH INTF block

The SEG_SCH INTF block receives data from the data block and outputs the data through a serial interface to the scheduler block. When the segmentation engine receives a status cell from the incoming SPI-4 interface, the data block outputs the data into the SEG_SCH INTF block instead of the data SRAM. As soon as the complete cell is

received, the SEG_SCH INTF block shifts the cell out the serial interface to the scheduler.

ERROR CONDITIONS

Several errors can be detected in the segmentation engine, and the error packet will be marked error by the SEG_ENGINE. The different errors are in Figure 143. These errors are reported in the status register. There are individual bits in the register to mask the individual errors from generating interrupt to the CPU.

Memory Description

Summary of Internal Memories:

See Figure 144.

Data structure for SEG_TABLE

See Figure 145. The SEG_TABLE is a 65 x 39 bits internal SRAM addressed by the port ID (PID). The segmentation table has 65 entries, one for each port and CPU.

BAD	:	This bit is set if detects error condition and needs to discard
SOP	:	This bit is set if start of packet is received for the port
EOP	:	This bit is set if end of packet is received for the port
Cell Byte Count	:	This field keeps track of the byte count received for the 64-byte cell
PTR	:	The start address to the cell in the data SRAM
Packet Length	:	This field keeps track of the total length for the incoming packet for the port
CLP	:	Cell Loss Priority, valid only with ATM traffic
EFCI	:	EFCI, valid only with ATM traffic
CALC_CRC	:	Calculate L2 CRC.
OAM	:	indicate OAM cell.

Data structure for FREE_BUFFER

See Figure 146. The FREE_BUFFER is an 80 x 7 bits internal SRAM that used to store the pointer to the free data block in the data SRAM. The segmentation engine takes a pointer from the top of this table to receive the incoming data. When the memory manager gets a complete cell off the queue the pointer to that cell is return to the end of this table.

Data structure for Q_FIFO

See Figure 147. This 80 x 70 bits internal SRAM that stores the queue of complete cells and the information for each queued cell.

FID	:	Flow ID of the incoming data for the port.
TYPE	:	Traffic type of the incoming data for the port.
BAD	:	Mark the packet to be discarded.
SOP	:	Mark the start of the packet.
EOP	:	Mark the end of the packet.
CELL_LEN	:	Indicate the number of valid bytes in the 64-byte cell.
PTR	:	The start address to the cell in the data SRAM.

- 1 PKT_LEN : Length of the packet, valid only for packet traffic.
- 2 QOS : Quality of Service.
- 3 O_PORT : Output port number.
- 4 CLP : Cell Loss Priority, valid only with ATM traffic.
- 5 EFCI : EFCI, valid only with ATM traffic.
- 6 CALC_CRC : Calculate L2 CRC.
- 7 OAM : indicate OAM cell.

8

9 **Data structure for DATA SRAM**

- 10 See Figure 148. The DATA SRAM is to store the incoming data in 64-bytes
- 11 block.

12

13 **Data structure for CRC_TABLE**

- 14 See Figure 149. The CRC Table is a 65 x 32 bits dual-port internal SRAM
- 15 addressed by the port ID (PID). Each port has an entry to this SRAM to store the partial
- 16 CRC generated for the AAL-5 traffic.

17

18 **Data structure for STATISTIC**

- 19 See Figure 150. The Statistic SRAM is a 65 x 80 bits internal SRAM addressed
- 20 by the port ID (PID). Each port has an entry to this SRAM to store the number of packets
- 21 and data received for that port. Statistics for the CPU interface is in the last entry. The
- 22 CPU interface accesses the content of this SRAM.

- 23 PKT_RCVD : Number of packets received for each port
- 24 DATA_RCVD : Number of data bytes received for each port

25

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

Interface

Interface Ports

System

See Figure 151.

Segmentation ⇔ Lookup

See Figure 152.

Segmentation ⇔ Memory Manager

See Figure 153.

Segmentation ⇔ PFQ

See Figure 154.

Segmentation ⇔ SCH

See Figure 155.

Segmentation ⇔ CPU Interface

See Figure 156.

Timing Diagrams

See Figure 157.

Memory Access

See Figure 158. The bolded “**Rd**” indicates that there are data dependencies for these reads. The current PID being processed is compared with the incoming one, the segmentation uses the current per port parameters if they are equal. Otherwise, it uses those read from the memories.

1

2 **Registers Description**

3

4 **Registers**

5 See Figure 159.

6

7 **Commands**

8

8 **Summary of commands**

9

See Figure 160.

10

10 **CPU Inject Packet**

11

See Figure 161.

12

12 **Read SEG FBL memory**

13

See Figure 162.

14

14 **Write SEG FBL memory**

15

See Figure 163.

16

16 **Read SEG CRC memory**

17

See Figure 164.

18

18 **Write SEG CRC memory**

19

See Figure 165.

20

20 **Read SEG STATISTIC memory**

21

See Figure 166.

22

22 **Write SEG STATISTIC memory**

23

See Figure 167.

24

24 **Read SEG QUEUE memory**

25

See Figure 168.

- 1 **Write SEG QUEUE memory**
- 2 See Figure 169.
- 3 **Read SEG DATA memory**
- 4 See Figure 170.
- 5 **Write SEG DATA memory**
- 6 See Figure 171.
- 7 **Read SEG TBL memory**
- 8 See Figure 172.
- 9 **Write SEG TBL memory**
- 10 See Figure 173.

Principles of Operation

The segmentation engine has three different sources for its incoming data: SPI-4, CPU interface, and status cells from per flow Q. The control information for the complete cells are queue into the queue FIFO, and are output to the memory manager following global sync scheme. If there is a status cell from per flow Q, it will be output in the next available cell slot, and the CPU slot will be used to compensate for the port that its slot is taken by the status cell. The segmentation engine resumes to the global sync in the following global sync cycle.

In normal operation, the segmentation gets its data from the SPI-4. If the FIFO full signal is not asserted, it means the segmentation engine is ready and has room in its FIFO to receive data from the SPI-4. The look up engine then sends the control information together with the data to the segmentation engine. For the first data cell of a packet for a particular port the segmentation engine sets the SOP bit for that port in the segmentation table, writes data into the data FIFO, increments both the packet length and cell count as data come in. The data state machine keeps receiving data from the look up engine until the end of data burst, complete cell size, end of packet, or if there is data coming from the CPU or serial interface.

If it is the end of data burst only, the segmentation engine loads the packet length and cell count back into the segmentation table for that port. The segmentation engine is now ready to receive the next data burst.

If the segmentation engine has a complete cell size, it queues the cell by writing the cell information into the Q_FIFO. The segmentation engine resets the cell count, loads the packet length back into the segmentation table. Again, the segmentation engine is now ready to receive the next data burst.

If it is the end of a packet, the segmentation engine checks to see if the traffic type is for packet to AAL-5. If so, the segmentation engine checks to see if another new cell is needed for the AAL-5 trailer. The segmentation engine then queues the completed cell(s) into the Q_FIFO, and then resets the entire entry for that port in the segmentation table. For packet data to AAL-5 cells, the CRC_GEN calculates the ATM CRC-32 when the cells are going out to the memory manager. The partial CRC is stored on per port basis, and the partial CRC is loaded back into the CRC generation engine for each cell queued out the particular port.

The segmentation engine can process multiple open ports. For each port ID it reads from the segmentation table the parameters for that port, receives data for that port, stores the parameters back into the table, read the parameters for a different port from the table and so on.

The segmentation engine queues a completed cell into the Q_FIFO and asserts the SEG_MEM_AVAILABLE signal to the memory manager block. The memory manager block pops the Q_FIFO to get the information (FID, PID, etc.), and the segmentation engine reads the data from the data SRAM using the pointer from the Q_FIFO. The segmentation engine sends the data to the memory manager with timing according to the global sync. After the segmentation engine sends the completed cell to the memory manager, it returns the pointer to the FBL.

If the incoming data is from the CPU interface, the segmentation engine receives data from the CPU_FIFO, strips the header from the cell and extracts the control

1 information from the header. The segmentation engine then waits to complete the cell and
2 queues the completed cell to the Q_FIFO as in normal operation.

3 If the incoming data is from the serial interface from per flow Q, the segmentation
4 shifts the serial data in, packs to 64-bit bus and store in an internal FIFO. The header of
5 the incoming cell is stripped and control information such as flow ID and traffic type are
6 extracted from the header. The segmentation engine waits to complete the 64-byte cell
7 and queues it into the control FIFO (Q_FIFO). When the segmentation engine detects a
8 status cell as the traffic type coming through the SPI interface it queues the complete cell
9 into a separate FIFO in the SEG_SCH_INTF block and send the data out in a serial
10 format to the scheduler.

11 **Memory Manager**

12 **Features**

13
14 The Memory Manager (MM) block manages all the data and control traffics
15 from/to the Segmentation, Per Flow Queue, Reassembly, and external SDRAM (or
16 DRAM). Input cells (data) are sent from the Segmentation block to the MM block. The
17 ATM cells and packets are sent to the Memory Manager block in the "internal cell"
18 format. Internal cell size is determined at the power CPU of the chip (i.e. programmable).
19 Figure 175 depicts the top level interface of the Internal Memory Manager block with
20 others.

21
22 The Memory Manager block has several main features listed as follows:

- 23 1. Support up to 512 Mbytes of external storage Memory using DDR-
24 SDRAM.
- 25 2. Cache the incoming "bypass" cells so that they can be transferred to the
26 Reassembly block with minimum latency (without going through the
27 external SDRAM path).
- 28 3. Act as a buffer for data going from the Segmentation block to the external
29 SDRAM.
- 30 4. Act as a buffer between the incoming burst from the external SDRAM and
31 the output to the Reassembly block.
- 32 5. Receive Cell Data and Control Information from the Segmentation block
33 and initiate En-queuing process.
- 34 6. Receive De-queuing information (such as BID, FID) from the Per Flow
35 Queue, get requested data from either the Internal Memory or the external
36 SDRAM, and then send them to the Reassembly block along with other
37 related information.
- 38 7. Interface with the Free Buffer List (FBL) inside the Per Flow Queue
39 block:
 - 40 a. Send En-queuing request to FBL to receive a BID indicating the
41 En-queuing cell's data location in DDR-SDRAM.
 - 42 b. Send De-queuing BID to FBL to release the location pointed by
43 the BID.
- 44 8. Manage the processes and orders of reading from the SDRAM and writing
45 to the SDRAM to maximize the overall bandwidth of the SDRAM. This

relates to the strategy of selection of banks of memory to read and write to reduce the overhead penalty.

9. CPU Interface: All memories of the Memory Manager can be accessed (read/write) by the microprocessor through the CPU interface:
 - a. Internal Memory
 - b. CAM, Internal Memory List (IML), Free Internal Memory List (FIML)
 - c. external SDRAM (by Queuing CPU read/write requests)
10. Frequencies: 200 Mhz

See Figure 174.

Functional Description

Refer to Figure 174 for the top level block diagram of the Memory Manager (MM) block and Figure 175 for the block diagram of the Internal Memory Manager (IMM). The incoming cells are sent to the Internal Memory Manager block by the Segmentation block along with other information for the En-queuing process. After getting the En-queuing BID from the Free Buffer List block, the IMM synchronizes the BID with all the information from the Segmentation block and then send them to the Per Flow Queue block for En-queuing.

Incoming cells are temporarily stored in the Internal Memory (internal SRAM) that is controlled by the Internal Write Manager block. If cells are specified as "bypass" indicated by the Type input, they stay in the Internal Memory and then they are sent to the Reassembly block directly during the De-queuing process. If incoming cells are specified as "non-bypass" indicated by the Type input, they are temporarily stored in the Internal Memory, and then moved to the external SDRAM. The moving of all "non-bypass" data is handled by processing all the SDRAM write requests with En-queuing BID/IMBID inside the SDRAM Write Request Queues block. During the De-queuing process, if some of these cells are called to be de-queued (i.e. CAM Hit) before they are sent to the SDRAM, their status will be changed in the CAM and they will stay in the Internal Memory instead of moving out to the SDRAM. These cells will be sent to the Reassembly block directly. Right after a cell is taken out of the Internal Memory, the CAM and the Free Internal Buffer List are Updated immediately.

During the De-queuing process, the Per Flow Queue interfaces with the Internal Memory Manager block with a De-queuing BID. The IMM will check in the CAM to see whether the requested cell (data) is inside the Internal Memory (SRAM) or inside the external SDRAM. If CAM indicates CAM Hit and the requested data is a "bypass" type, it means that the requested cell is inside the Internal Memory. In that case, the Memory ID Manager (MIM) will issue a command to the Internal Read Manager to read the CAM hit data out of the Internal Memory; the IMBID (internal buffer ID) will be provided to the Internal Read Manager. The MIM also modifies the CAM table so that this specific data will not go to the external SDRAM (it will be marked as "bypass") and it signals the De-Queue Request Manager block not to send a read SDRAM request to the SDRAM Read Request Queues block.

During the De-queuing process, if the De-queuing BID can not be found inside the CAM, it means that the requested data is inside the external SDRAM. In this case MIM will signal the De-Queue Request Manager block to send a read SDRAM request to

the SDRAM Read Request Queues block with the BID and IMBID (provided by the MIM block from the Free Internal Buffer List). The SDRAM Read Manager block takes read requests from the SDRAM Read Request Queues and processes them to through the SDRAM Manager (outside of the Internal Memory Manager block). Read-back data from the SDRAM will go through the SDRAM Read Manager block with the data itself, the BID and the IMBID assigned to it previously. The data and IMBID will go to the Internal Write Manager which will open the gate for the data to be stored inside the Internal Memory block at the IMBID location. As data enters into the Internal Memory, the CAM table will also be Updated properly (it will be marked as "bypass"). The Internal Read Manager then will schedule the data to go out to the Reassembly properly based on the order of the De-queue requests came in previously (stored in the De-queue Queue block).

Microprocessor can access (read/write) the external SDRAM through the CPU data bus for testing and diagnostic purposes. The access can be random and data are not managed or shaped. The CPU Interface block inside the Memory Manager sends read or write requests from/to the SDRAM Read Request Queues and the SDRAM Write Request Queues, respectively. The CPUID/BID will be queued along with the IMBID/BID. The SDRAM Read Manager and the SDRAM Write Manager will send these information and associated data to the SDRAM BW Optimizer based on the bank selection strategy inside it. Data coming back from the SDRAM is routed properly to the CPU Interface block based on the returned CPUBID/BID information.

Block Description

Internal Memory Manager

Refer to Figure 175 for the block diagram of the Internal Memory Manager.

Internal Memory Block

This is an internal dual-port SRAM memory. It can run at 200 Mhz with 5 ns access time. The input data bus width is 128 bits and the output data bus width is also 128 bits. The input data, IMBID and control signals are provided by the Internal Write Manager block. The output data go to the Internal Read Manager in response to the read request along with IMBID provided by the Internal Read Manager block.

The size of the SRAM is currently defined as 64 Kbytes (= 4K x 128 = 1K of 64-byte cells).

Memory ID Manager (MIM) Block

This block consists of a CAM, a Free Internal Memory List (FIML), an Internal Memory List (IML), and a MIM (Memory ID Manager) Controller. Refer to Figure 3 and Figure 4 to understand the mechanism of the CAM, FIML, and the IML blocks. FIML and IML are internal SRAMs. The CAM contains the BID values. The MIM Controller block manages all the incoming commands to Update the CAM, FIML, and IML blocks.

1 It also arbitrates all the Update commands from different blocks as they can happen at the
2 same time.

4 **Free Internal Memory List (FIML):**

5 This list contains 1K entries (i.e. to support 1024 cells inside the Internal Memory
6 (SSRAM)). The content of each entry is the IMBID that is an address points to the
7 Internal Memory List (IML) to get IMBID and status. It also corresponds to the output
8 address of the CAM. The FIML contents are initialized from 0 to 1023 when an
9 initialization command is received from the CPU interface.

10 The FIML has 2 pointers: TakePtr and ReturnPtr. TakePtr and ReturnPtr are
11 initialized at the top of the List, ReturnPtr points to location 0 and TakePtr points to
12 location 1. TakePtr points to a free IMBID value. TakePtr is incremented by 1 every time
13 an IMBID is taken. ReturnPtr points to a not-available IMBID value. Every time a
14 location is returned to the FIML, a new (i.e. returned) IMBID value will be written into
15 the location that ReturnPtr currently points to. Now that location represents a free cell
16 with the returned IMBID points to a free cell inside the Internal Memory. ReturnPtr is
17 also incremented by 1 every time an IMBID is returned.

18 When TakePtr and ReturnPtr reach above the value of 1023, they will be rolled
19 over to the value of 00. TakePtr and ReturnPtr will never be equal. If (TakePtr –
20 ReturnPtr) is equal to -1, then the IM is full and the IM_Full bit of the
21 Mem_Status_0_reg will be set. If (TakePtr – ReturnPtr) is equal to 1, then the IM is
22 empty and the IM_Empty bit of the Mem_Status_0_reg will be set. The status of Full-1 is
23 also recorded in the IM_Full_1 bit.

24 The CPU interface should be able to write and read to/from the FIML memory:

25 1). CPU_Wr_FIML{input: Adr, Wr_value}:

26 This command allows the CPU interface to write to the FIML table where Adr is the
27 address value (range from 0 to 1023) and Wr_value is the content that we are going to
28 write to. Refer to Register Description section.

29 2). CPU_Rd_FIML{input: Adr; output: Rd_value}:

30 This command will read the content of the FIML table at location Adr. Refer to Register
31 Description section.

33 **Internal Memory List (IML):**

34 The IMBID value is the address of the IML block. The content of each entry
35 contains 3 bits: the valid bit, bypass bit, and readback bit.

- 36 ■ Valid: this bit indicates that at IMBID location in the Internal Memory, a
37 cell is valid.
- 38 ■ Bypass: IMBID points to a cell that is whether a bypass or non-bypass
39 type. A bypass type cell will not go into the external SDRAM, it is ready
40 inside the IM to go out to the Reassembly. A non-bypass type cell is in
41 transition to go into the external SDRAM.
- 42 ■ Readback: indicates that this cell is from the SDRAM and ready to go out
43 to the Reassembly block.

45 **CAM:**

CAM contains the values of BIDs (external SDRAM pointers) that are inside the Internal Memory. If a CAM hit occurs, an IMBID value will be returned by CAM. After receiving the IMBID, the MIM will go to the IML and read the content at IMBID. This content contains IMBID (Internal Memory address) and its status. If valid bit is equal to 1, it means that IMBID is valid and, therefore, this CAM hit is valid. Otherwise, a “miss” is returned.

There are 4 main internal commands (between blocks and not CPU commands) can be issued to the MIM block (from other blocks):

1). Take_ID \leftarrow IMBID:

In this command, the MIM will access the FIML and take a free IMBID from it. The output of the command will be a free IMBID that represents a free cell inside the Internal Memory. The TakePtr will be incremented by 1 after it releases a free cell.

2). Return_ID{input: IMBID}:

In this command, the MIM will receive an IMBID as an input. It will return this pointer to the FIML. The ReturnPtr will be incremented by 1 after an IMBID is returned.

3). Compare_CAM{input: BID; output: IMBID, status, hit/miss}:

In this command, input argument is BID and output is IMBID, status (valid bit, bypass bit, busy bit) and hit/miss bit. The MIM receives BID and compares it to the entries in the table. If a match is found, the CAM will output the corresponding IMBID. The MIM goes to the IML and read its content at the location IMBID. The MIM also checks to see if the “valid” bit is equal to 1 at this location. If it is equal to 0 then a miss has occurred, otherwise, it is a CAM hit. The content of IML at IMBID will include status, and a hit/miss bit.

4). Wr_CAM{input: IMBID, BID, status}:

This command will modify the content of the CAM table and the content of the IML table with the new values of BID and status, respectively.

5). Rd_CAM{input: IMBID; output: BID, status}

This command reads the contents of the CAM table and the IML table.

The CPU interface should be able to write and read to/from the CAM table:

1). CPU_Wr_CAM{input: IMBID, BID, status}:

This command allows the CPU interface to write to the CAM table and the IML table with new values of BID, and status, respectively.

2). CPU_Rd_CAM{input: IMBID; output: BID, status}:

This command will read the content of CAM table (i.e. BID) and the content of IML table (i.e. status)

Internal Write Manager Block

This block receives input from the Segmentation block and the SDRAM Read Manager, and manages the writing process into the Internal Memory (SRAM) block.

There are 4 main functions that this block handles independently:

Function 1: (Segmentation \rightarrow Internal Memory).

- 1) Segmentation asserts Seg_Mem_Available signal to the MM block to let it know that there is data available in its fifo.
- 2) MM block sends Mem_Seg_Cntl_Pop signal to pop the cell control information of a cell.
- 3) Cell control information is sent to the MM block.
- 4) At a fixed delay from the time that IWM gets the cell control information from the Segmentation block, the Segmentation block starts to send cell data to the MM block at 128-bit word at 10ns period (every other cycle of 200 Mhz clock). The 128-bit data word will be first latched into flip-flops inside the IWM.
- 5) From FIML, get free IMBID: issue Take_ID command to MIM block.
- 6) Write data into the Internal Memory: Wr_IM{IMBID; Seg_Data}.
- 7) Update CAM with new BID: Wr_CAM{IMBID; BID, valid=1, bypass=0/1, readback=0}.
- 8) If (not "bypass"),
 - a. Then put new BID/IMBID into end of SDRAM Write Request Queues: Wr_SWRQ{end of queue; BID/IMBID}.

Note: This block has a small FIFO to receive data inputs from the Segmentation block. The size of the FIFO is 4 64-byte cells with 128 bit wide. Data will be stored in this IWM FIFO before it is sent to the Internal Memory.

Function 2: (Segmentation → Per Flow Queue).

- 1) IWM sends the cell control information to the PFQ block with a command specified in Mem_PFQ_Com signals (i.e. "00": invalid; "01": enqueueing; "10": discard; "11": release (for dequeuing)).
- 2) PFQ receives the information, processes it and then
 - a. If command = "01", then it sends back en-queue BID (Buffer ID) to the MM block 8 cycles later. This BID will be used as an address to write data into the external SDRAM. If BID = "7FFFFFFF", then that particular cell will be dropped and not written into the SDRAM.
 - b. If command = "10", then it sends BID = "7FFFFFFF" back to the MM block. In this case, data will not be written into the external SDRAM.

Function 3: (Control data path to Internal Memory).

Inputs and outputs of the Internal Memory are time shared as follows (to handle dual-port 128-bit wide SRAM at 200 Mhz):

- 1) Slot 0 (period = 5ns): input from Segmentation.
- 2) Slot 1 (period = 5ns): input from SDRAM Read Manager.

Function 4: Microprocessor (CPU) Access.

When switched to this mode, the microprocessor can write to the Internal Memory (IM) with a specified address. During this mode, there is only the CPU interface block that can access (read/write) to the IM block. Refer to the Register Description section..

CPU_Wr_IM{input: Adr, Wr_value}: write to IM at the Adr location with the value of Wr_value.

1

2 **SDRAM Read Manager Block**

3 This block receives read requests from the SDRAM Read Request Queues and
 4 transfers those requests with data to the SDRAM Manager. It receives read-back data
 5 from the SDRAM Manager and sends it to the Internal Memory.

6

7 **Function:** (SDRAM → Internal Memory).

- 8 1) SDRAM Read Manager sends data with BID and IMBID to the Internal Write
 9 Manager: $Rd_SRRQ\{top\ of\ queue\} \leq BID, IMBID, SDRAM_Data$.
- 10 2) Request the Internal Write Manager block to write data into Internal Memory:
 11 $Wr_IM\{IMBID; SDRAM_data\}$.
- 12 3) Update CAM with new BID (send a command to the MIM): $Wr_CAM\{IMBID; BID,$
 13 $valid=1, bypass=1, readback=1\}$.

14

15 **Internal Read Manager Block**

16 This block receives input from the Internal Memory and manages its output to the
 17 Reassembly block and the SDRAM Write Manager.

18

19 **Function 1:** (Internal Memory → Reassembly).

- 20 1) Look into CAM to see if BID from top of De-queue Queue matches with BID in
 21 CAM. Issue a command to MIM: $Compare_CAM\{BID; IMBID, status, hit/miss\}$. If
 22 $hit/miss = 1$, get IMBID as output from CAM. Put IMBID into the output queue
 23 inside the IRM block. Assert $Mem_Ras_Available$ to the Reassembly block.
- 24 2) The Reassembly block asserts the $Ras_Mem_Cntl_Pop$ signal to the MM block. The
 25 IRM sends the cell control information at the top of the queue by popping from top of
 26 the De-queue Queue block.
- 27 3) At a specified time slot the Reassembly block asserts the $Ras_Mem_Data_Pop$ signal
 28 to the MM block. The IRM starts to read data from the IM, $Rd_IM\{IMBID\} \leq$
 29 $Reassembly_Data$.
- 30 4) Send data to the Reassembly block. Pop IMBID from the output queue.
- 31 5) If $(REL=1)$ or $(REL=0)$ and " $bypass$ " = 1 and " $readback$ " = 1)
 32 Then return IMBID to Free Internal Memory List: Return $ID\{IMBID\}$.
- 33 6) Update CAM: $Wr_CAM\{IMBID; BID=00, valid=0, bypass=0, readback=0\}$.
- 34 7) If $(REL = 1)$ then
 35 Release BID to FBL in PFQ: assert Mem_PFQ_Valid high, $Mem_PFQ_Com =$
 36 "11" (release BID) and $Mem_PFQ_Param = BID$.
 37 (Refer to Multicast and BID/IMBID Release section)

38

39 **Function 2:** (Control data path to Reassembly and SDRAM).

40 Inputs and outputs of the Internal Memory are time shared as follows (to handle a dual-
 41 port 128-bit wide SRAM at 200 Mhz):

42

- 1 1. Slot 1 (period = 5ns): output to Reassembly.
- 2 2. Slot 2 (period = 5ns): output to SDRAM Write Manager.

4 **Function 3: Microprocessor Access.**

5 When switched to this mode, the microprocessor can read from the Internal Memory (IM)
6 at a specified address. During this mode, there is only the CPU interface block that can
7 access (read/write) to the IM block.

8 CPU_Rd_IM{input: Adr; output: Rd_value}: read from IM at the Adr location with the
9 value of Rd_value.

11 **SDRAM Write Manager Block**

12 This block receives write requests from the SRAM Write Request Queues and
13 then gets data from the Internal Memory to send to the SDRAM Manager block.

15 **Operation: (Internal Memory → SDRAM).**

- 16 1) Get BID and IMBID from the top of SDRAM Write Request Queues (refer to
17 SDRAM read/write strategy section for more details on algorithm to select proper
18 banks of memory to read and to write).
- 19 2) Look into MIM to check status of that BID: Rd_CAM{IMBID; BID, status}.
- 20 3) If "bypass" = 1,
21 Then do not write data to the SDRAM. Return BID to FBL. Go to last
22 step (5).
23 Else if "bypass" = 0, then get data from the Internal Memory: Rd_IM{IMBID} <=
24 IM_Data. Send IM_Data to the SDRAM Manager along with BID and IMBID.
- 25 4) Return IMBID to the Free Internal Memory List and Update CAM:
26 Return_ID{IMBID}
27 Wr_CAM{IMBID; BID=00, valid=0, bypass=0, readback=0}.
- 28 5) Pop the BID and IMBID from the top of SDRAM Write Request Queues.

30 **De-queue Request Manager Block**

31 This block receives De-queueing requests from the Per Flow Queue block with
32 other related information. It will look into the CAM to see if there is a CAM hit. If it is, it
33 will modify CAM to reflect its new status, that is now a "bypass" data and should not be
34 written into the SDRAM. It then lets the Internal Read Manager take data out of the
35 Internal Memory and send it to the Reassembly block. If there is no CAM hit, it will send
36 a read request to the SDRAM Read Request Queues to get data out from the SDRAM.

38 **Function: (Per Flow Queue → Internal Memory, SDRAM).**

- 39 a. Store BID, PortID and other cell control related information for De-
40 queueing process into the end of De-queue Queue.
- 41 b. Look into CAM to see if BID is matched with BID in the CAM.
42 Compare_CAM{BID; IMBID, status, hit/miss}

- 1 c. If BID is found in CAM (hit/miss = 1) and "readback" = 0 (not outgoing
 2 cell)
 3 If "bypass" = 0
 4 Then get IMBID. Modify CAM: Wr_CAM{IMBID; BID, valid=1,
 5 bypass=1, readback=0}.
 6 Else if no CAM hit,
 7 Then get free IMBID from Free Internal Memory List: Take_ID. Write
 8 BID and IMBID into end of SDRAM Read Request Queue.
 9

10 Multicast and BID/IMBID Release

11 During Dequeueing the PFQ sends to the MM block a command to indicate
 12 whether to release the Dequeue BID to the Free Buffer List (in the PFQ) or not to release
 13 the Dequeue BID due to Multicast. Depending on the type of data inside the Internal
 14 Memory, the IMBID is also released to the Free Internal Buffer List (FIBL) in the MIM
 15 block. Figures 176-178 depicts the mechanism of the release command:

16 SDRAM Manager (SM)

17 See Figure 179. The DDR_SDRAM Controller (SM) receives the memory write
 18 (enqueue) and read (dequeue) requests from the Internal Memory Manager (IMM). All
 19 operations are done on a per controller basis independent of the other controller. Each of
 20 the two SDRAM Manager blocks, SM0 and SM1 consists of:

- 21 • RD ID FIFO
- 22 • WR ID FIFO
- 23 • BANDWIDTH OPTIMIZER BLOCK
- 24 • REQUEST SERVICE BLOCK
- 25 • CONTROLLER INTERFACE BLOCK
- 26 • BUFID ADDR DECODE BLOCK
- 27 • DENALI DDR_SDRAM MEMORY CONTROLLER BLOCK

28 SM0 and SM1 blocks each control corresponding pair of SDRAMs. It uses 128-bit
 29 data bus at 200 MHz in order to support 2X OC-192 transfer rate (20 Gbps). Each
 30 DDR_SDRAM Controller has the capability of peak data transfers at 16 bytes (or 128
 31 bits) per 5ns or 25.6 Gbps for a total memory bandwidth of 51.2 Gbps in the
 32 DDR_SDRAM Manager (SM). Therefore, the DDR_SDRAM Manager must
 33 simultaneously process two memory accesses, one each to SM0 and SM1 in order to
 34 sustain the required bandwidth.

35 The **BANDWIDTH OPTIMIZER** block decides which of 8 possible operations,
 36 4 reads from each bank and 4 writes to each bank to service for optimal data throughput
 37 from the SDRAMs. It looks at the current state which is the last processed request and
 38 the outstanding requests in the per bank queues and picks one based on an algorithm
 39 detailed below to optimize the data throughput from/to the SDRAMs. Some of the
 40 avoided situations include sequential same bank operations with their bank activate and
 41 precharge delays, read followed by write commands with the associated pitfall of turn
 42 around cycle delays, etc.
 43

The Bandwidth Optimizer Block will scan through the 8 queues, 4 reads and 4 writes per Memory Manager, DM0 and DM1, and select an operation request to the appropriate bank based on the SDRAM Read/Write Strategy discussed below.

- Each series of Read requests to external memory will have to follow by a Write request to a different bank of memory from the accessed bank of the last Read.
- Within a series of Read requests, consecutive Reads should have different banks.
- Within a series of Write requests, consecutive Writes should have different banks.

For Example:

$Rd(b0) \rightarrow Rd(b1) \rightarrow Rd(b2) \rightarrow Rd(b3) \rightarrow Wr(b0) \rightarrow Wr(b1) \rightarrow Wr(b3) \dots$

At the same time the read and write request BID/IMBID are stored in the **READ ID FIFO** and **WRITE ID FIFO** respectively to be returned to the IMM with the read data or a write request confirmation respectively. For write request, IMBID is returned so that IMM can update the CAM and release IMBID to Free Internal Memory List. For read request, BID/IMBID are returned so that IMM can store the read data to IM and release BID to Free Buffer List.

The **REQUEST SERVICE BLOCK** is activated with the particular request to be serviced by the Bandwidth Optimizer block. It routes the data back and forth appropriately to and from the dram and communicates with the Address Decode block and the Memory Controller Interface block. It also handles the back pressure from the Denali Memory Controller in case the command queues or the data queues in that block have reached their limit.

The **BUFFERID ADDRESS DECODE BLOCK** does the address translation of the BID to the absolute byte address and feeds this to the Memory Controller Interface block via the Request Service Block.

The **CONTROLLER INTERFACE BLOCK** formats all the commands for the Denali Memory Controller. It processes the address and data part and sends them to the Memory Controller along with required control signals in the format expected by the Memory Controller. It also receives the data from the memory controller for read operations and sends it to the Request Service Block. It will pass on control signals from the Controller to the Request Service block signaling that the controller command or data queue is full and when it is ready to receive the next request

From the above block onwards the interface of the **DENALI MEMORY CONTROLLER** and the SDRAMs is detailed in the diagram below. The Memory Controller takes in the Address and Data from the Memory Controller Interface Block and formats them into appropriate RAS and CAS commands for the DRAM memories. It also does the refresh of the SDRAM memories.

See Figure 180.

CPU Interface

Assumption: CPU access requests are infrequent and are for debugging and testing purpose only.

SM services the CPU request along with the en-queuing and de-queuing requests with equal priority. Whenever CPU requests an access to DDR-SDRAMs, SM will process the request at the end of one round scanning and servicing the 8 queues of the read/write requests. In order not to interrupt the bandwidth for enqueueing and dequeuing,

each CPU access will be serviced after X number of rounds of servicing the read and write queues. (The X number is programmable through a register)

Since the en-queuing cells and de-queuing cells are 64-byte and hence the DDR SDRAMs are programmed to have burst length of 4 ($64 \times 8 = 4 \times 128$), the CPU data access should be in the same cell unit, 64-byte.

Memory Description

Summary of Memories:

See Figure 181.

Summary of CAM:

See Figure 182.

Interface

Segmentation Engine

See Figure 183.

Per Flow Queue Engine

See Figure 184.

DDR-SDRAM

See Figure 185.

Reassembly

See Figure 186.

CPU Interface

See Figure 187.

Global Sync

See Figure 188.

Interface with DDR-SDRAM

The Memory Manager can support up to 512 Mbytes of external memories using DDR-SDRAM. The DDR-SDRAM device that is available at this time is the 1M x 32bit x 4banks (128 Mbits). The memory size can be programmed using register Mem_Config[2:0] (refer to the Register section).

2 DDR-SDRAM devices are required to support a data bus of 64bit. Hence 2 data strobes, Mem_SDR_DQS_Lo and Mem_SDR_DQS_Hi, are used to for the lower and upper 32bit of the data bus respectively.

I A duplication set (set 1) of memory address and control signals (including memory clock) is available for each controller in order to improve its signal driving capability when the number of devices are increasing (up to 16 devices per controller). Hence each interfacing signals will not drive more than 8 devices in any configuration of memory sizes.

This second set of address and control pins can be activated by programming the register Mem_Config[3].

Example 1:

Refer to Figure 189 for the configuration of 64 Mbytes of memory.
Mem_Config[3] = 0.

Example 2:

Refer to Figure 190 for the configuration of 256 Mbytes of memory.
Mem_Config[3] = 1.

Register Description

See Figure 191.

Commands

CPU_Wr_SDRAM

See Figure 192.

CPU_Rd_SDRAM

See Figure 193.

CPU_Wr_FIML

See Figure 194.

CPU_Rd_FIML

See Figure 195.

CPU_Wr_CAM

See Figure 196.

CPU_Rd_CAM

See Figure 197.

CPU_Init_MM

See Figure 198.

Principle of Operation

1 Initialization Setup

2 At power-on-reset, the entire block is reset by MEM_Reset_L signal. After
3 power-on-reset, all registers are at default values.

4 CPU can issue the CPU_Init_MM command (refer to the "Register section") that
5 will go through proper steps to initialize all internal memories inside the Memory
6 Manager block.

8 Operational Flows

10 Normal Mode:

11 The normal mode is the default mode. The Memory Manager handles both the
12 enqueueing and the dequeuing processes.

14 *Enqueue:*

16 Segmentation to MM:

- 17 1. Segmentation block has data available for MM block, it starts to assert
18 Seg_Mem_Available signal.
- 19 2. MM block asserts Mem_Seg_Cntl_Pop signal to get cell control information from
20 Segmentation block.
- 21 3. At a specified time slot the Segmentation block sends cell control information to the
22 MM block.
- 23 4. The MM block stores cell control information in the IWM block.
- 24 5. At a specified time slot the Segmentation block starts to send data to the MM block
25 through the Data bus (128 bits) in 4 clock cycles at 10 ns period (every other clock
26 tick of the 200 Mhz clock).
- 27 6. The IWM block stores data inside the IWM FIFO (4 64-byte cells).
- 28 7. After receiving the free BID from PFQ (see below), the IWM block write appropriate
29 data from the IWM FIFO into the Internal Memory (IM).

32 MM to PFQ:

33 The IWM block sends cell control information to the PFQ block on the
34 Mem_PFQ_Param along with a command via the Mem_PFQ_Com bus. The
35 Mem_PFQ_Valid signal is asserted high to let the PFQ block knows that the data bus and
36 command bus contains valid data. If the Discard signal is not asserted then the
37 Mem_PFQ_Com value should be 01 (enqueue). . If the Discard signal is asserted then the
38 Mem_PFQ_Com value should be 10 (discard).

- 39 1. If it is an enqueue command, the PFQ block sends back to the MM block a BID at a
40 specified number of clocks later via the PFQ_Mem_Param bus. This BID is used as an
41 address to write data into the external SDRAM. The PFQ_Mem_Valid signal is
42 asserted and the PFQ_Mem_Com value is 01.
- 43 2. IWM block receives BID and tags it along with the associated data being stored in the
44 IWM FIFO.

3. If it is a discard command, the PFQ block sends back a BID with value of "7FFFFFFF" to the MM block at that specified time slot. In this case, data from the Write FIFO will not be written inside the external SDRAM (it will be flushed out of the FIFO).
4. If the PFQ decides to drop the cell (even when the command is not a discard), it will send to the MM block a BID with value of "7FFFFFFF". In this case, data from the Write FIFO will not be written inside the external SDRAM (it will be flushed out of the FIFO).

MM to SDRAM:

1. Data from the Segmentation is temporarily stored in the IWM FIFO.
2. After receiving a free and valid BID from the PFQ, the IWM block requests a free IMBID from the Memory ID Manager (MIM) block. The CAM is also updated properly with the new IMBID and BID.
3. IWM writes data from the IWM FIFO to the Internal Memory (IM) block at address IMBID.
4. If the cell is not a "bypass" cell, then IWM sends BID and IMBID to the SDRAM Write Request Queue (SWRQ) block. Place the information at the bottom of the queue that associates with memory bank of the BID address.
5. The SDRAM Write Manager (SWM) block goes through the "Write Bank Selection Algorithm" and it selects a BID/IMBID pair from a particular queue inside the SWRQ block. The SWRQ sends that IMBID to the MIM block for status verification.
6. If that cell is marked as a "bypass" cell (this means that it is being dequeued and ready to go out to the Reassembly block), then do not write data to the SDRAM. Pop the BID/IMBID from the top of the queue inside the SWRQ block.
7. If that cell is not valid as indicated in the status (this means that it has gone out to the Reassembly block), then do not write data to the SDRAM. Pop the BID/IMBID from the top of the queue inside the SWRQ block.
8. If that cell is marked as a "non-bypass" cell, then
 - a). The SWM block sends request to the IRM block with BID/IMBID.
 - b). IRM reads data from the IM block (Rd_IM{IMBID}) and send this data to the SDRAM Manager block along with BID/IMBID.
 - c). Return IMBID to the MIM block and update its CAM:
 - Return_ID{IMBID}
 - Wr_CAM{IMBID; BID=00, valid=0, bypass=0, readback=0}
 - Pop BID/IMBID from the top of the queue inside the SWRQ block

Dequeue:

PFQ to MM:

1. At a specified time slot, PFQ sends to the Read Request Queue & Manager (RRQM) block dequeued cell control information that includes a BID via the Data2 bus. Both Valid2 and Deq signals are asserted to let the MM block knows that the data bus is valid and contains dequeue information.
2. The RRQM block stores cell control information at the end of the queue.
3. The RRQM block sends a BID from the top of the queue with a read request to the Memory Read/Write Controller (MRWC).

SRAM to MM:

1. The MRWC block monitors the Read FIFO. If it is not full, the MRWC sends a dequeue BID to the SRAM Interface (SI) block with a read command ("10").
2. The SI block puts the BID and the command in a queue.
3. The SI block executes a read command to the external SRAM with the BID as an address. Data is read out from the external SSRAM and sends to the Read FIFO. The associated BID will be placed inside the Read FIFO along with the data.
4. When a complete cell is read back from the SSRAM, the Read FIFO sends its status and the associated BID to the MRWC block.
5. The MRWC block sends the BID back to the PFQ block via the MM_PFQ_Data bus. The Valid is asserted and the Com value is "11" (dequeue).
6. PFQ receives the BID and releases it back to the Free Buffer List block.

MM to Reassembly:

1. The Read Request Queue & Manager (RRQM) block monitors the Read FIFO status.
2. Whenever there is a complete cell available inside the Read FIFO, the RRQM sends the MM_Ras_Available to the Reassembly block.
3. The Reassembly block asserts the Cntl_Pop signal to the MM block. The RRQM block sends the cell control information at the top of the queue to the Reassembly block at a specified time slot.
4. At a specified time slot the Reassembly block asserts the Data_Pop signal to the MM block. The Read FIFO starts to send data out to the Reassembly block at a specified time slot via the 128-bit Data bus at 10 ns period (every other clock tick of the 200 Mhz clock).

Flow charts for Functional Operations

See Figures 199-205.

Reassembly

The Reassembly Engine is a state machine that reassembles the stored switch cells into outgoing ATM, AAL-5, packets, or switch cells.

Main Features

- Can provide a maximum of 1M total flows.
- Uses PER PORT queues for re-assembly.
- Reassembles segmented cells back into IP packets. These packets were segmented into cells for switching at the ingress side and need to be reassembled back
- Maximum packet size of 64 KB can be handled
- AAL5 trailer field check for length and CRC-32 for all flows that have to be re-assembled and sent out as packets.
- L2 CRC check for packets that have the L2 CRC check enabled.
- Can strip up to 16 bytes of Header data.
- Reassembles 48 byte AAL5 cells or 52 byte ATM cells, both stored as 64 byte memory cells into outgoing packets, ATM cells or switch cells.
- No lookup, pass-through traffic
- For ATM
 - o VPI, VCI header translation
 - o Operation, Administration and Maintenance (OAM), Cell Loss Priority (CLP) and Explicit Forward Congestion Indication (EFCI) replacement
 - o Mark End of Packet in Payload Type Indicator (PTI) Field for AAL5 Egress Cells
 - o Optional MPLS tag append
 - o Encapsulate raw cells into packets for routing through packet fabrics.
 - o Encapsulate ATM cells into packets for routing through packet fabric.
- MPLS lookup
 - o POS (PPP/IP) packet format
 - o MPLS push/pop
 - o MPLS header translation
 - o Decrement TTL and drop if TTL is one before decrementing
- It has one re-assembly going on for each of the ports that are conFigured and active, supports a maximum of 64 ports plus a CPU port.
- The output of this block goes to the PHY PORT.

- Supports the traffic represented in Figure 206.
- Applications for various traffic types
- Types of reassembly:
 - o TYPE 1: Ingress TM (ATM, AAL5, Packets) to 64B Switch Cells
 - o TYPE 2: Ingress ATM cells to Packet Data
 - o TYPE 3: Ingress AAL5 cells to Packet Data
 - o TYPE 4: Ingress Packets to Packet Data
 - o TYPE 5: Egress (ATM) to ATM cells with Translation
 - o TYPE 6: Egress (AAL5) to Packet Data
 - o TYPE 7: Egress (Packets) to AAL5 Cells
 - o TYPE 8: Egress (Packets) to Packet Data
- Calculates and checks the CRC-32 for AAL-5 reassembly packets
- Each port has a 32-bit counter to count packets/cells transmitting from that port
- Each port has a 48-bit counter to count number of data bytes transmitted through that port
- Each port has a 16-bit counter to count packet CRC errors.
- Each port has a 16-bit counter to count TTL Timeout errors.
- Internal memory to store control and data for the 64 ports
 - o Dual-port SRAM
 - CPU Data - 32 x 64 bits to store CPU data
 - Data - 1024 x 64 bits to store Cell data
 - Header - 256 x 64 bits to store Header data
 - Output Port Control - 128 x 21 bits to store Output Port Control information and Next Pointer links.
 - Output Port Queue - 64 x 24 bits to store Output Port Queue Pointers to the Data cells
 - CRC - 64 x 32 bits to store the partial CRC for each port
 - o Single-port SRAM
 - Reassembly Table - 64 x 18 bits to store packet length, SOP and EOP
 - Free Buffer - 128 x 10 bits to maintain the free buffer list for data SRAM
 - Statistic - 64 x 96 bits to store packets, data and errors received for each port
- Operates at 200 MHz

Functional Description

See Figure 207. The Reassembly Engine reassembles the outgoing data from cells stored internally into ATM cells, Switch cells, or Packet data based on the application.

As shown in Figure 206 there could be 8 different types of reassembly. The operation of the Reassembly Engine for each of these types is explained in detail in the following sections.

1 TYPE 1 Reassembly - Ingress Memory cells => Switch cells

2 Type 1 reassembly maps the 64-byte cells in Memory into switch cells for
3 transmission by adding a Switch Header. This Reassembly type handles Application
4 types 0-3 & 7. Figure 2 shows the one to one mapping of ATM Memory cells into switch
5 cells. The switch header will be 8 or 16 bytes long with the required cell fabric header
6 being first and our 4-byte header following with the remaining bytes being Pad. Type 1
7 Reassembly also informs the SPI interface of the number of valid bytes in each 8-byte
8 word and which Cell is the SOP and which Cell is the EOP for Packet traffic. For AAL5
9 cells the cells were was stored in memory as AAL5 switch cells with the AAL5 trailer at
10 the end and Pad to fill the cell to 64 Bytes. For Packet traffic the cells were stored in
11 memory as switch cells with an AAL5 like trailer added to pass on the packet length
12 information. Figure 3 shows the AAL5 encapsulated cells and Figure 4 shows the pseudo
13 AAL5 encapsulated packets. For Ingress TM types 0 & 1 that transfer ATM cells the
14 OAM, EFCI and CLP bits will be inserted in our 4 byte switch header but for Ingress TM
15 packets, types 2 & 3, the bits will be set to zero.
16 See Figures 208-10.

18 TYPE 2 Reassembly - 52 Byte ATM cells => Packet Data

19 Type 2 Reassembly maps the 52-byte ATM cells stored as 64 byte cells in
20 Memory into packets and adds a packet header before sending it on to the SPI interface.
21 Only one ATM cell will be mapped into each packet. Eight of the 12 bytes of Pad data at
22 the end of each cell will be removed before the packet is sent out. This will allow us to
23 send out a 64-byte packet including 8-bytes of header. If the Switch Packet Header is 16-
24 bytes then the packet will be 72-bytes long. The OAM, EFCI and CLP bits will be
25 passed in our 4 bytes of the Switch Packet Header.

26 See Figure 211.

TYPE 3 Reassembly - 48 Byte ATM cells => Packet Data

Type 3 Reassembly maps the 48-byte ATM cells stored in Memory as 64 byte cells into packets by stripping the 16 bytes of PAD from each cell and then optionally stripping the L2 Packet Header and adding a Switch Packet Header before sending it on to the SPI interface. It must check the L2 CRC to ensure the packets integrity before sending on to the switch fabric as it also checks the CRC in the AAL5 Trailer and removes the AAL5 encapsulation. The OAM, EFCI and CLP bits will be passed in our 4 bytes of the Switch Packet Header. It must check the L2 CRC to ensure the packets integrity before sending it to the switch fabric as the L2 CRC being sent out is now incorrect due to the removal of the L2 Header and must be replaced by the Egress chip when it adds a new L2 Packet Header. The Lookup engine will replace the L2 Header and the Segmentation will calculate the new CRC.

See Figure 212.

TYPE 4 Reassembly - Switch cells => Packet Data

Type 4 Reassembly maps the AAL5 like Switch cells stored in memory into packets removes the AAL5 encapsulation and may add an 8 or 16-byte Switch Packet header. This can also operate without a lookup operation. (Traffic Shaper, Packet in, Packet out)

See Figure 213.

TYPE 5 Reassembly - 52 Byte ATM cells + PAD => ATM cells with Translation

Type 5 Reassembly translates the ATM Header by replacing the 28-bit VPI/VCID then replaces the OAM, CLP & EFCI bits before the cell is sent out. The extra 12 Bytes of Pad data from the end of each ATM cell stored as 64-byte cells in Memory is stripped before sending it out as a 52-byte ATM cell. Using the data byte enables of the SPI interface strips the extra Pad.

See Figure 214.

TYPE 6 Reassembly - 48 Byte ATM cells => Packet Data

Type 6 Reassembly maps the 48-byte AAL5 encapsulated ATM cells stored in Memory as 64-byte cells into packets by stripping the 16 bytes of PAD from each cell and then performing a MPLS operation based on the contents of the Flow ID look-up table before sending it on to the SPI interface. It also checks the CRC in the AAL5 Trailer before removing the AAL5 encapsulation. If the MPLS TTL is one before being decremented the packet will be dropped in Reassembly.

See Figure 215.

TYPE 7 Reassembly - 48 Byte ATM cells => AAL5 cells

Type 7 reassembly adds the 4 byte ATM Header along with 4 bytes of PAD, so it is 8 byte word aligned, to the 48-byte AAL5 encapsulated ATM cells stored in Memory as 64-byte cells, strips the 16 bytes of Pad, replaces the OAM, CLP & EFCI bits in the ATM Header and then transmits them as AAL-5 cells. The 4 bytes of alignment PAD will be stripped using the byte enables when transmitting to the SPI. The EOP AAL-5 cell is marked using the PTI field in the ATM Header.

See Figure 216.

TYPE 8 Reassembly - Switch cells => Packet Data

Type 8 Reassembly maps the AAL5 like Switch cells stored in memory into packets and performs a MPLS operation based on the contents of the Flow ID look-up table. It also checks the CRC in the AAL5 Trailer before removing the AAL5 encapsulation. If the MPLS TTL is one before being decremented the packet will be dropped in Reassembly.

See Figure 217.

Block Descriptions

See Figure 218.

Description of INPUT/ENQUEUE_SM Block

The CTRL_SM block receives PID, FID, TYPE, EOP, SOP, CLP, EFCI and MEM_RAS_AVAILABLE from the Memory Manager. When AVAILABLE is asserted the CTRL_SM pops the control information from the Memory Manager control FIFO. The state machine then reads from the control SRAM using PID as the address to retrieve parameters for that port. The control state machine also requests a pointer to a free buffer from the Free Buffer List (FBL). Then it loads the parameters into the internal registers and starts the DATA_SM. The type of Reassembly is decoded from the TYPE field by the control state machine and the enable signal for each type is passed to the data state machine. The SOP and EOP signals inform the Reassembly engine when the first and last cells of a packet are received.

The FID is used to read the external Header Look-Up SRAM to retrieve the appropriate header and MPLS information for that Flow Id. This Header is passed on to the DATA_SM to be added to the reassembled cell.

The PID and other control signals are read from the Memory Manager by the Reassembly engine 1 clock before the data to give the control state machine enough time to read the per port parameters from the SRAM and load into internal registers, retrieve a free buffer pointer from the FBL and to start the header read from external memory.

Description of OUTPUT/DEQUEUE_SM Block

The data state machine starts when it receives the start signal from the control state machine. Once started, the data state machine reads 64-bytes of data from the Memory Manager data FIFO. The data state machine will also determine from the TYPE field whether to check the CRC.

When a complete cell is read and the data state machine has added the appropriate header information obtained from the external Look-Up SRAM, the data state machine informs the control state machine that it is finished reading the data for that cell. The control state machine then stores the per-port parameters back in the reassembly control table. The data state machine is then ready to receive more data.

“Once the next cell has been received and determined not to be an EOP cell then the previous cell for that port is queued into the Q_FIFO for transmission to the SPI. If the next cell is an EOP then the whole cell must be read and the length determined from

the AAL5 trailer to ensure the proper amount of valid data in the previous cell before it is queued for transmission.

Once the data for a cell has been stored and the control information is queued in the Q_FIFO, the Reassembly asserts READY to the SPI interface so it can come and Pop the Q_FIFO for the control information for that cell. Once the Q_FIFO has been Popped by the SPI, the Reassembly will start reading the data associated with that control information and send it on to the SPI 64 bits per 200 Mhz clock cycle until the complete cell is transmitted. A Data_Valid signal will be active when the data is valid on the bus and 8 byte enable signals will inform the SPI which bytes have valid data in each 64 bit transfer. This operation will allow the SPI to control the flow of data by controlling when the Q_FIFO is Popped and also strip the extra pad bytes when necessary.

The data state machine has two types of counters, one for operation and one for statistic. The cell byte count and packet length counters are for operation. The packet received count and data received count and CRC errors are for statistic. All counters are per port basis.

Description of CPU I/F Block

This is a standardized interface to the CPU in the system. It will be instantiated within each block that requires read/write access to configuration registers or internal and external memory arrays. During system operation the reassembly engine will use this interface for writing configuration registers and the header lookup table in external SSRAM as well as to receive status and statistic information. The CPU will also receive Data Cells or Packets through this interface using the CPU Port Buffer (port 65).

Description of CPU_PORT Buffer Block

This is a 32 x 64 bit internal Dual Port SRAM to store the incoming data on a per memory cell basis such that the Reassembly Engine can reassemble the data into the proper cell configuration before sending it to the CPU. The memory is organized as 4 64-byte buffers to allow 4 cells for the CPU port to be stored.

Description of REASM_TABLE_SRAM Block

This is a 64 x 18 bit dual-port internal SRAM to store the parameters for the 64 ports. This RAM is accessed twice (1 read and 1 write) by the control state machine for each data cell that is received for a particular port.

Description of FBL Block

See Figure 219. This is a 128 x 10 bit single-port internal SRAM to store the free pointers to the cells in the data SRAM. This list needs to be initialized by software after reset to have the free pointers. The control state machine requests a free pointer by reading the list, and the data state machine returns a pointer by writing back to the end of the list when the SPI receives a complete cell from the Reassembly block.

Description of CRC_GEN Block

This block generates the CRC that covers the PDU (Protocol Data Unit) in the AAL-5 format. The polynomial used is CRC-32 as shown in Figure 220. There is a minimum of 3 clocks delay in generating the CRC because the internal data path is 64

bits and the CRC polynomial is calculated 32 bits at a time. The CRC is generated in 2 stages, first the upper 32 bits and then the lower 32 bits.

The reassembly engine calculates CRC-32 for the AAL-5 data. The CRC_GEN block reads the partial CRC with port ID as the address to the CRC SRAM. The CRC_GEN block loads the partial CRC into the CRC engine and starts calculating as the data is coming in from the memory manager. When the complete cell is loaded in the internal SRAM the CRC_GEN block writes back the partial CRC into the CRC SRAM, again using port ID as the address. At the end of a packet, the CRC is compared to the CRC in the AAL-5 trailer and the packet is marked as bad when sent to the SPI out interface if the compare fails.

Description of CRC_TBL Block

This is a 64 x 32 bit dual-port internal SRAM to store the partial CRC for each port. The data state machine controls the read or write to this table based on the traffic type.

Description of DATA_PIPE Block

The DATA_PIPE block delays the data to the Data SRAM so that it matches the delay to lookup and inserts the appropriate header information and generates and check the CRC-32 for all Packet reassembly.

Description of DATA_SRAM Block

This is a 1024 x 64 bit internal Dual Port SRAM to store the incoming data on a per memory cell basis such that the Reassembly Engine can reassemble the data into the proper cell configuration before sending it to the SPI. The memory is organized as 128 64-byte buffers to allow 2 cells for each port to be stored.

Description of HEADER_SRAM Block

This is a 256 x 64 bit internal Dual Port SRAM to store the Header data on a per memory cell basis such that the Reassembly Engine can reassemble the data into the proper cell configuration and attach a programmable Header to it before sending it to the SPI. The memory is organized as 128 16-byte buffers to allow 2 16-byte Headers for each port to be stored.

Description of OUTPUT PORT QUEUE Block

This is a 64 x 30 bit single-port internal SRAM to queue information of the complete cells. The data state machine queues the cell control information when a complete cell is ready to send to the SPI interface. The Reassembly Engine sends an Available signal to the SPIO when the Q_FIFO is not empty telling the SPIO that at least one cell of data is available for transmission.

Control information delineates the SOP, EOP and SOB as well as the Port Number and whether the packet data had a CRC error or another error condition occurred. The FIFO also holds the 10-bit pointer to the data cell so the Reassembly will send the associated data with that control information.

Description of OUTPUT CONTROL_SRAM Block

This is a 64 x 30 bit single-port internal SRAM to queue information of the complete cells. The data state machine queues the cell control information when a complete cell is ready to send to the SPI interface. The Reassembly Engine sends an Available signal to the SPIO when the Q_FIFO is not empty telling the SPIO that at least one cell of data is available for transmission.

Control information delineates the SOP, EOP and SOB as well as the Port Number and whether the packet data had a CRC error or another error condition occurred. The FIFO also holds the 10-bit pointer to the data cell so the Reassembly will send the associated data with that control information.

Description of External Look-Up Memory

This is a 2M x 72 bit external SRAM to store the Header control information per FLOWID. This includes the appropriately programmed new header for the particular FLOW or the necessary MPLS manipulation control information for that FLOW. This can handle 1M FLOWS with 16 Bytes of header and 2 bytes of control information per FLOW.

Description of MPLS Label Manipulation

The FLOW-ID Lookup table will contain the MPLS label operation instructions and new labels. Three bits in the first control word will contain the encoded instructions for what MPLS operation to perform. The operations will be either replace the MPLS tag, Push the MPLS tag once or twice or Pop the MPLS tag once or twice. After each MPLS operation the original MPLS TTL value will be decremented. If the TTL value is one before being decremented then the packet must be dropped. This will be done by not enqueueing the packet into the output Queue FIFO and subsequently dropping the following cells from the FLOW until the EOP is reached.

If two labels are being pushed then Label 1 will be pushed on first and Label 2 will be pushed on second. If one label is pushed label 1 will be used. If two labels are popped the third label will be the one remaining on the top and if one label is popped then the second label will be remaining on top.

The TTL will be decremented as follows:

- Replace MPLS Tag – Decrement MPLS TTL
- Push 1 MPLS Tag – Decrement original MPLS TTL
- Push 2 MPLS Tags – Decrement original MPLS TTL
- Pop 1 MPLS Tag – Decrement original MPLS TTL
- Pop 2 MPLS Tags – Decrement original MPLS TTL

TAGS

Figures 221 and 222 show the data formats for each of the supported protocols. The specifications are not hard coded. The port table will be used for the header's location. Each field is 8 bits wide with 12 bits VPI and 16 bits VCI field, for a total of 28 bits lookup tag. This 28 bit lookup tag is located in the 1st 64-bit word, bits 32 thru 63. The PPP protocol is part of the link layer header. The tag is bit 60:41 of the first 64 bits word.

1 ATM TYPES

2 If the type of the traffic is ATM, then the CLP, EFCI, and OAM bits are included
3 in the Switch Header for Ingress operation and inserted into the ATM Header for Egress
4 operation. See Figure 223.

6 Description of Memories

7 There is one external Memory interface and 10 internal memory blocks for the
8 Reassembly Engine. There is a 64 x 18 bit single port SRAM for a Reassembly Table
9 SRAM, a 128 x 18 bit dual-port SRAM for Output Port Control information, a 64 x 18 bit
10 dual-port SRAM for Output Port Queue Pointer information, a 128 x 7 bit single port
11 SRAM for a Free Buffer list, a 1024 x 64 bit dual-port SRAM for DATA SRAM, a 256 x
12 64 bit dual-port SRAM for HEADER DATA, a 36 x 64 bit dual-port SRAM for CPU
13 DATA, two 64 x 32 bit dual-port SRAMs for CRC Tables and a 64 x 112 bit single port
14 SRAM for Statistic information.

16 Data Structure for Reassembly Table

17 See Figure 224. The REASSEMBLY_TABLE is a 64 x 18 bit internal SRAM
18 addressed by the port ID (PID). The Reassembly Table has 64 entries, one for each port.

- 19 ○ SOP: Start of Packet
- 20 ○ EOP: End of Packet
- 21 ○ Packet Length: A running count of the number of bytes in that packet.

24 Data Structure for OUTPUT_CONTROL_MEMORY

25 See Figure 225. The OUTPUT_CONTROL_MEMORY is a 128 x 18 bit internal
26 SRAM that is used to store the control information for each cell that is ready to send to
27 the SPI block along with the next read pointer and next number of Header words for the
28 next cell linked to that port in the DATA SRAM and the number of valid bytes in that
29 cell.

- 30 ○ SOP: Start of Packet
- 31 ○ EOP: End of Packet
- 32 ○ CELL_LENGTH: Number of Data bytes valid for that cell. Used
33 internally to create the Byte valid signals to the SPI interface.
34 (Combine with number of Header words for
35 RAS_SPI_CELL_LENGTH sent to SPI OUT.)
- 36 ○ NEXT_RD_PTR: Pointer to the Head of the next cell of data queued
37 to that port in the Data SRAM.
- 38 ○ NEXT_HDR_WD: Number of header bytes for the next cell of data
39 queued to that port in the Data SRAM.

41 Data Structure for OUTPUT_PORT_QUEUE

42 See Figure 226. The OUTPUT_PORT_QUEUE is a 64 x 18 bit internal SRAM
43 that is used to store the output port Queue Head (RD_PTR) and Queue Tail (WR_PTR)
44 pointers for each cell that is ready to send to the SPI block along with the port Empty
45 status and the number of Header words in the Queue Head cell.

- RD_PTR: Pointer to the Head of the next cell for that port that should be sent to the SPI Out.
- WR_PTR: Pointer to the Head of the last cell for that port that should be sent to the SPI Out.
- E: Empty – Status of cells queued and linked for that port. (Default = 1, Empty)
- PEND: Cell Pending – Cell is waiting to be queued to output port. (Default = 0, Not Pending)
- HDR_WDS: Number of header words for the next cell for that port to be sent to the SPI Out.
-

Data Structure for FREE_BUFFER

See Figure 227. The FREE_BUFFER is a 128 x 7 bit internal SRAM that is used to store the pointer to the free data block in the DATA SRAM. A pointer is taken from the top of this table when a data block is needed to receive the incoming data. When the SPI interface gets a complete cell off the queue the pointer to that cell is returned to the end of this table.

- FREE_PTR: Pointer to the head of a 64-byte buffer in the Data SRAM.

Data Structure for DATA_SRAM

See Figure 228. The DATA_SRAM is used to store the data while reassembling it into the proper cells for transmission. There are 128 64-byte buffers. This can support 128 64-byte cells.

Data Structure for HEADER_SRAM

See Figure 229. The DATA_SRAM is used to store the data while reassembling it into the proper cells for transmission. There are 128 16-byte buffers. This can support 128 16-byte Headers.

Data Structure for CPU_DATA_SRAM

See Figure 230. The CPU_DATA_SRAM is used to store cell data destined for the CPU. There are 4 72-byte buffers. This can support 4 64-byte Cells of Data plus 8 bytes of control information.

Data Structure for CRC_TABLE

See Figure 231. The CRC_TABLE is a 64 x 32 bit internal dual port SRAM that stores the partial CRC for each port reassembly.

- PART_CRC: Partial CRC stored after each cell for that port. Used to calculate the CRC for an entire packet.

Data Structure for STATISTICS

See Figure 232. The Statistic SRAM is a 64 x 112 bit internal SRAM addressed by the port ID (PID). Each port has an entry to this SRAM to store the number of packets and data received for that port along with the number of CRC error packets. The content of this SRAM is accessed through the PCI interface block.

- PKT_RCVD: Counter for the number of Packets received.
- DATA_RCVD: Counter for the number of Data bytes received. Only includes valid data bytes.
- CRC_ERRORS: Counter for the number of CRC Errors received.
- TTL_ERRORS: Counter for the number of TTL Timeout Errors.

Data Structure for EXTERNAL LOOK-UP MEMORY

See Figure 233. The External Look-Up SRAM is 4M x 36 bits made up of eight 512K x 36 bit external SRAMs addressed by the FLOW ID (FID). Each FLOW has two entries to this SRAM to store the Header data for that Flow. The Flow Header Data (FHD) is two or four 36-bit words that contain up to 2 MPLS tags, 1 VPI/VCI label, or a Layer 2 Header. On the Ingress chip a total of 12 bytes are available for new user defined headers with the remaining 4 bytes of Header data used for a Maximus internally defined Switch Header. On the Egress chip all 16 bytes are available for new user defined headers as our Maximus 4 byte header will not be append in this case. The two control bytes contain a number of fields that indicate how to process the packet or cell. Each FHD is linked to a user defined Header, an MPLS tag or a VPI/VCI label through the FLOWID.

FLOW HEADER FORMAT:

See Figure 234. The following fields are in the Flow Header Data Table:

1. 1 Byte Control word in 1st FLOWID word to determines the type of FHD (MPLS, ATM or L2) and the MPLS operation to perform as well as the number of Header words.
2. 32 ATM/32 MPLS bit New Label 1
3. 32-bit MPLS New Label 2
4. 1 Byte Control word in 2nd FLOWID word to determine the number of L2 Header bytes to be stripped.
5. 32 bits for possible Longer Header Formats.
6. 32 bits for our internal Switch Header information.

CONTROL BYTE FORMAT:

See Figure 235. The first control word is bits 71 to 64 of the first 72-bit word of the FLOW ID with the second control word being bits 71 to 64 of the second 72-bit word for that FLOW ID. The following fields are in the Flow Header Control Words:

1. HDR TYPE = 2 bits that determines the type of Header (MPLS, ATM or Ethernet)
 - a. 00 = MPLS
 - b. 10 = ATM
 - c. 01 = L2
 - d. 11 = Reserved
2. MPLS OPERATION = 3 bits to indicate the type of MPLS processing to perform.
 - e. 000 = Replace MPLS Tag; Decrement MPLS TTL
 - f. 001 = Push 1 MPLS Tag; Decrement original MPLS TTL
 - g. 010 = Push 2 MPLS Tags; Decrement original MPLS TTL

- h. 011 = Pop 1 MPLS Tag; Decrement original MPLS TTL
- i. 100 = Pop 2 MPLS Tags; Decrement original MPLS TTL
- j. The last 3 cases are reserved for future implementations.
- 3. HDR = 2 bits to indicate the number of valid Header words. (00 = 0, 01=1, 1X=2)
- 4. L2EN = L2 CRC checking Enabled
- 5. # BYTES L2 HDR STRIP = 4 bits for Number of L2 Header Bytes to strip if type is ETHERNET.
- 6. RSVD = 4 bits Reserved for future use.

MAXIMUS 4 BYTE HEADER FORMAT:

See Figure 236. The Ingress Maximus chip will append a proprietary 4-byte header onto the end of any 8 or 16 byte switch or packet header before sending the data out. This allows the user to append a programmable 4 to 12 byte header along with the Maximus 4-byte header onto every packet or switch cell that is transmitted from the Ingress Maximus chip. The Maximus header contains information of Flow ID, application types and Class, along with the OAM, CLP and EFCI bits for ATM cells.

- o FID: FLOW ID
- o TYPE: Application Type
- o EFCI: Explicit Forward Congestion Indication
- o CLP: Cell Loss Priority
- o OAM: Operation, Administration and Maintenance
- o CLASS: Class of service
- o EOP: End Of Packet
- o SOP: Start Of Packet

The FID and CLASS will be obtained from the Header Lookup memory while the remaining signals will be obtained from the Memory Manager interface.

ATM HEADER FORMAT:

See Figure 237.

- o VPI: Virtual Path Identifier
- o VCI: Virtual Channel Identifier
- o PT: Payload Type
 - EOP: End of Packet
 - EFCI: Explicit Forward Congestion Indication
 - MI: Management Information
- o CLP: Cell Loss Priority
- o HEC: Header Error Control

MPLS HEADER FORMAT:

See Figure 238.

AAL5 TRAILER FORMAT:

See Figure 239.

- o CPCS-UU: Common Part Convergence Sublayer
- o CPI: Common Part Indicator

- L: Length of payload
- CRC: Error detection in SSCS PDU (user data)

Interface

System

See Figure 240.

Reassembly ⇔ Memory Manager

See Figure 241.

Reassembly ⇔ SPI

See Figure 242.

Reassembly ⇔ CPU

See Figure 243.

TIMING

See Figures 244-6.

Memories

Input phase to RAS from Internal Memory Manager. See Figure 247. Note: R2 for the Partial CRC Table is dependent on the PID. If the next cell has the same PID as the previous cell then this read information is not used.

ERROR CONDITIONS

Several errors can happen in the Reassembly engine with the error packets being marked during the EOP transfer to the SPI as EOP Abort. The different errors are listed in Figure 248.

Register Description

Registers

See Figure 249.

Commands:

See Figure 250.

Commands

Read external RAS memory

See Figure 251.

Write external RAS memory

See Figure 252.

1 FLOW HEADER FORMAT:

2 See Figure 253. The following fields are in the Flow Header Data Table:

- 3 7. 1 Byte Control word in 1st FLOWID word to determines the type of FHD
4 (MPLS, ATM or Ethernet) and the MPLS operation to perform as well as the
5 number of Header words.
6 8. 32 ATM/32 MPLS bit New Label 1
7 9. 32-bit MPLS New Label 2
8 10. 1 Byte Control word in 2nd FLOWID word to determine the number of
9 Ethernet L2 Header bytes to be stripped.
10 11. 32 bits for possible Longer Header Formats.
11 12. 32 bits for our internal Switch Header information.
12

13 Read Internal Data Memory

14 See Figure 254.

16 Write Internal Data Memory

17 See Figure 255.

**19 Read Four Internal Memories (Partial CRC1, Partial CRC2, Reassembly Table,
20 Output Queue)**

21 See Figure 256.
22

**23 Write Four Internal Memories (Partial CRC1, Partial CRC2, Reassembly Table,
24 Output Queue)**

25 See Figure 257.
26

27 Read PORT Statistics

28 See Figure 258.
29

30 Write PORT Statistics

31 See Figure 259.

1

2 **Read CPU Data Memory**

3 See Figure 260.

4

5 **Write CPU Data Memory**

6 See Figure 261.

7

8 **Read 2 Internal Memories (Free Buffer List (FBL), Output Control)**

9 See Figure 262.

10

11 **Write 2 Internal Memories (Free Buffer List (FBL), Output Control)**

12 See Figure 263.

13

14 **Read Internal Header Data Memory**

15 See Figure 264.

16

17 **Write Internal Header Data Memory**

18 See Figure 265.

19

20 **Get CPU Cell command**

21 See Figure 266. This command reads 64 bytes from the CPU Data FIFO 8 bytes
22 at a time and stores them in the general purpose registers R0- R15. There is no address
23 associated with this command as it just reads the CPU Data FIFO 8 times and then loads
24 the Registers with the data.

1

2 **Initialize Free Buffer List and Output Port Queue Memories command**

3 See Figure 267. This command initializes the Free Buffer List in the FBL
4 Memory and sets the Empty and Pending bits to inactive in the Output Port Queue
5 Memory.

6

7 **Reserved Op-codes**

8 Op-codes 17 to 31 are reserved for future use.

9

10 **Principle of Operation**

11 The re-assembly block keeps track of a FLOW ID per port. The scheduler has a
12 weighted round robin scheme that it uses to schedule a CELL for each port. If a port has
13 more weight (or bandwidth) then it schedules more cells for that port. These scheduled
14 cells are provided to the de-queue engine which sends the requests on to the Internal
15 Memory Manager for pulling them from the SRAM.

16 The Internal Memory Manager reads the scheduled cell from the external SRAM
17 and stores it in an Output Data FIFO. Once the cell is stored in the Output Data FIFO the
18 Internal Memory Manager queues the control information for that cell into the Output
19 Control FIFO and asserts an available signal to the Reassembly Engine.

20 The Reassembly Engine pops the Output Control FIFO to get the control
21 information for the available cell. It uses the FLOW ID to look up the Header processing
22 necessary for that cell as well as the Type field to determine the type of Reassembly to
23 perform.

24 It also obtains a pointer to the Internal Dual Port Data SRAM from the FBL FIFO.
25 Once it has started the Header look up and obtained a Free Buffer Pointer the Reassembly
26 Engine starts reading the data from the Internal Memory Manager Output Data FIFO.
27 The data is read 128 bits at a time at 100 MHz. Each Data cell is 64 bytes long, with
28 each data buffer being 64 bytes deep. There are 128 64-byte buffers in the Internal Data
29 SRAM. The Header information will be stored in a separate Internal Dual Port SRAM
30 that is also addressed from the FBL Pointer. The Header buffer has 128 16-byte buffers
31 so it can handle up to 16 bytes per Header. The Header will be stored concurrently with
32 the Cell Data being stored in the data buffer.

33 Once the Header and Data are stored locally in memory the cell control
34 information will be queued in the Output Control Memory and the Buffer pointer will be
35 queued in the Output Port Queue Memory. If more than one cell is queued on a port then
36 the Next Read Pointer will be stored in the Output Control Memory of the previous cell
37 to link the current cell to the previous cell.

38 The SPI has a weighted round robin scheme that it uses to Request a CELL for
39 each port. If a port has more weight (or bandwidth) then it Requests more cells for that
40 port. When the SPI requests a cell for a particular port the Reassembly checks the ports
41 Output Port Queue to see if there are any cells available for that port. If there are no cells

1 available the Reassembly informs the SPI that the port has no data by asserting an Empty
2 signal and the SPI continues on to the next port. If the port has cells available the
3 Reassembly uses the Read pointer stored in the Output Port Queue to read the control
4 information for that cell along with the associated Header and Cell Data. The
5 Reassembly sends the Header and Cell Data to the SPI interface along with a valid signal
6 and byte enables to specify the valid data bytes. The SPI will use the byte enables to strip
7 any unwanted padding. After the Reassembly Engine sends the complete cell to the SPI
8 interface it returns the Read Pointer to the FBL FIFO.

9 The Reassembly engine has a Reassembly status table for each of the ports. It can
10 simultaneously reassemble one packet at a time for each port. The local port Reassembly
11 status table stores the control information for each port along with state information and
12 or any partial cell information. The control information tells the re-assembly block how
13 to re-assemble the FLOW for that port, whether it is ATM to ATM, ATM to packet,
14 packet to ATM or packet to packet.

15 For AAL5 to Packet reassembly a CRC- 32 is calculated on the packet as each
16 cell is received. The partial CRC is stored on a per port basis and loaded back into the
17 CRC engine for each subsequent cell in that packet. Once it has received the end of the
18 packet the calculated CRC is check against the received CRC to ensure packet integrity.
19 If the CRC is in error the packet will be marked as bad by setting the Bad control bit in
20 the Output Queue FIFO for the SPI interface.

21 For Layer 2 Packet reassembly the CRC-32 is checked on the packet as it is being
22 sent to the SPI out interface. The CRC is calculated the same as the AAL5 CRC-32, but
23 if there is an error detected it must be sent immediately to the SPI out as the complete
24 packet has been sent in this case. The CRC checking is user selectable by programming
25 the CRC enable bit in the Flow ID Lookup table.

26 Statistics are kept on a per port basis for the number of Packets sent, the number
27 of Bytes sent and the number of errors received. The CPU can read and write internal
28 registers or statistics on a per port basis as well as access all internal memory locations
29 and the external Header Lookup memory through the PCI interface. Data packets can
30 also be sent to the CPU through the PCI interface.

31 32 **Initialization Setup**

33 The initialization process consists of writing the Initialize Memories command
34 op-code to the Command register. The Reassembly block will then initialize the Free
35 Buffer List and clear all the Empty and Pending bits in the Output Port Queue Memory.
36 When the initialization process is done, the command op-code will be cleared to all
37 zeroes to inform the system CPU that the initialization process is complete.

38 39 **Re-assembly flow**

- 40 1. Output scheduler selects the port.
- 41 2. Priority scheduler for port decides on the FLOW ID.
- 42 3. FLOW ID is sent to the De-queue engine.
- 43 4. Output scheduler keeps the current head and other information per port.
- 44 5. De-queue engine queues the buffer pointers along with corresponding port #
- 45 for Internal Memory Manager.

- 1 6. Internal Memory Manager gives a data available signal to the reassembly
- 2 engine when a complete cell is in the output data FIFO in the Internal Memory
- 3 Manager.
- 4 7. Internal Memory Manager sends ACK back to the scheduler
- 5 8. Scheduler releases the buffer from the INTERNAL MEMORY MANAGER
- 6 and passes next port and buffer head to the de-queue engine.
- 7 9. Internal Memory Manager tells Reassembly that a cell is available.
- 8 10. Reassembly engine pulls the control information for that cell from the Internal
- 9 Memory Manager control FIFO and then extracts the data from the Internal
- 10 Memory Manager output data FIFO.
- 11 11. Reassembly uses the control information to look up the Header processing
- 12 necessary for that FLOW. It also obtains a free buffer pointer to the Data
- 13 SRAM and uses the control information to determine the type of processing to
- 14 perform.
- 15 12. Reassembly engine stores the new Header information along with the data cell
- 16 and Queues the cell onto the appropriate output port.
- 17 13. The SPIO requests data on a per port basis based on its own internal port
- 18 calendar.
- 19 14. Once the SPIO requests data for a port that has a cell queued on that port the
- 20 Reassembly sends the control information for that cell along with the
- 21 associated Header and Cell Data for transmission.
- 22 15. The SPIO will strip extraneous pad data based on the byte enables that go with
- 23 the data.
- 24 16. The Reassembly engine sends CPU Buffer Full status information to the SPIO
- 25 for combining with the port FIFO status information that is sent to the
- 26 Scheduler from the SPIO.
- 27 .
- 28 .
- 29 .
- 30 .
- 31 .

Per Flow Queue Engine

The Per flow queue stores flows parameters and statistics in dedicated FID memory. In addition, it stores memory buffer pointers or IDs that hold data traffic of the flow IDs in a dedicated BID memory. Furthermore, it will retrieve the above FID parameters, statistics and memory buffer IDs when they are needed. It manages the free buffers list. It performs Random early discards on FID data cells.

Features

- Supports 1M possible flows
- Supports 1K Tunnel Flows.
- Supports 1K Multicast Flows.
- Multicast one flow to 1 K flows maximum.
- Collects statistics on a per flow basis.
 - Packet Counters: 32 Bits.
 - Byte Counters: 48 Bits.
 - Packet Discards per external Requests such as SPI: 16 Bits
 - Packet Discards per Reassembly Timeout request: 16 bits.
 - Packet Discards per internal buffer management request: 16 bits.
- Discards by returning buffer Ids to free buffer list.
- Manages the Free Buffer List.
- Manages RED by FID association.
- RED Memory: 1K x 48 bits.
- FID Enqueue Memory: 2M x 72 bits.
- FID Dequeue Memory: 1M x 36 bits.
- FID Statistics Memory: 2M x 72 bits.
- BIDLL Memory: 8M x 36 bits, including free buffer list.
- Tunneling Dequeue Leaf Memory: 1K x 56 bits, internal.
- Tunneling Enqueue Leaf Memory: 2K x 82 bits, internal.
- Tunneling Root Memory: 512 x 21 bits, internal.
- Multicast Dequeue Leaf Memory: 1K x 36 bits, internal.
- Multicast Enqueue Leaf Memory: 2K x 82 bits, internal.
- Multicast Root Memory: 512 x 30 bits, internal.
- Clock Frequency: 200 MHz.

Functional Description

See Figure 278. The Per Flow Queue Engine consists of the following blocks: Enqueue engine, Dequeue engine, Free List buffers manager, RED manager, Statistics engine, Timeout manager, and CIH Director. Enqueue engine places received data cell pointers on the queues. Dequeue engine retrieves and removes data cell pointers from the queues. The RED manager will track the total number of buffers on eight CLASS counters using each of the associated RED parameters. The cells will be dropped when the RED algorithms flags the cells as such. The statistics manager counts per flow various parameters. The timeout manager discards packets or cells when they exceed a certain timeout threshold on the enqueue side for the reassembly timeout. The timeout

manager discards packets or cells when they exceed a certain timeout threshold on the dequeue side for the reassembly timeout. The PFQ interfaces to the memory manager, to the shaper and to the scheduler. The Per Flow Queue Engine handles the flow tunneling, discarding and multicasting of FID in addition to enqueueing and dequeuing the FIDS. The statistics block will collect on a per flow basis the number of discarded flows, the number of bad flows, and the number of dropped flows because of RED. It will provide free buffers BIDS to the memory manager and it returns the buffer BIDS to the free buffer list.

Block Description

Algorithms

Enqueue Normal FID Algorithm

The enqueue request will place a BID of an FID on the queue. It will modify the BID as well as the FID Enqueue memory. When the queue is empty, the first enqueue is for an SOP cell. If the cell is not SOP, then the cells are dropped until the next SOP. See Figure 279.

1. The FID Enqueue memory is read, the tail pointer will point to the last BID in the link.
2. The BID tail is compared to "Null"; if it is a null then the queue is empty.
3. The current BID tail that was stored in the FID Enqueue memory BIDS will be used to access the BID memory. The new BID Tail is written as the content of that BIDS Tail location. If the cell is SOP or EOP then the flags are written in the BID memory with the BID link.
4. The new BID replaces the tail pointer in FID Enqueue memory.
5. The new tail is linked to "NULL" in the BID memory.
- When the tail pointer is read, it will be checked to see if it is null. When it is null then the Dequeue memory has to be written with the new head and the associated flags such as the SOP, CLP, and EOP etc.
- If an EOP is received, then a shape command is issued to the shaper.
- The count is sent to the statistics block with the FID when the EOP flag is received.
- The count is sent to the shaper and then it is not saved in the FID or BID memories.

Dequeue Normal FID Algorithm

The de-queue request will remove a BID from an FID.

1. The FID Dequeue memory is read; the head BID pointer will point to the first BID in the link.

2. The current BID head that was stored in the FID Dequeue memory will be used to access the BID memory.
3. The new BID will replace the head pointer and the value is written in the FID memory. If the BID has other flags like SOP, EOP they are written at the same time as the BID head.
- The new BID is compared to Null. If it is equal to null, then the Queue is empty and the Enqueue memory for that FID has to have the tail as Null.
- A dequeue command is sent to the memory manager.

See Figure 280.

Enqueue Multicast FID Algorithm

When the queue is empty, the first enqueue is for an SOP cell. If the cell is not SOP, then the cells are dropped until the next SOP. The bit in the control register will determine the treatment of multicast traffic. If the bit is set to treat the multicast traffic as Unicast, then the traffic is enqueued, as Unicast one. Else, the traffic is treated as described in the following section. The multicast control bit in the control register is reset. The PFQ will discard all the enqueue requests from the memory manager when an EOP cell is received for a multicast traffic. When the SOP cell is received, it is enqueued as a Unicast cell. The buffer ID pointers are linked in the memory. The pointers in memory are enqueued only once. The count is read from the Multicast Root memory. Every BID will have the count of the previous BID. When an SOP cell is enqueued, the tail pointer is in the enqueue memory. There is no link in the BID memory to the tail pointer in the enqueue FID yet. When a NULL is written to next location, the count is properly written in the correct field. When the EOP cell is received, then the BID link is performed as it was done for the previous cells. When the BID of the EOP cell is linked, the multicast enqueue operations has to link the head and tail BIDS to the various FIDS that are members of the multicast group. At this time, the enqueue engine will discard all the cells with the enqueue request from the memory manager until the multicast packet has been linked to all the FIDS in the multicast group. The multicast engine will perform the rest of the operations to enqueue the packet on the rest of the FIDS. The multicast engine will read the next FID that in the multicast group. The first FID member has already the pointers linked in its location. The multicast engine will update the head and the tail pointers in the enqueue and dequeue memories locations that correspond to the FID members if the FID queue is empty. If the FID is not empty, then the tail pointer is updated only. When the last FID member is update it, then the enqueue engine will continue to enqueue rather than discard the cells. The process is not efficient and the performance suffers considerably.

The buffer management whether it is RED or class will be performed on the first SOP when the cell is enqueued for the HEAD member of the multicast group. Thereafter there are no checks done. After receiving the EOP cell, the DBS will get the FIDS that are members of the multicast group only.

1. The packet is enqueued as Unicast one until the EOP cell.
2. When EOP cell is received, then the packet has to be linked for the rest of the FID members of the multicast group.
3. The FID root memory is read to retrieve the head FID, tail FID pointers.

4. The head FID has the packet already linked on its queue. The next FID field in the multicast group is read.
5. The next FID is compared to the tail FID to determine the last FID in the multicast group.
6. The LEAF memory of the FID in the next FID field is read to determine if the queue is empty. If the queue of the FID is not empty, then the tail pointer is updated it with the last pointer of the EOP cell. The new tail pointer is linked to the old tail pointer in the BID memory and the COUNT is written there, too.
7. If it is empty, then the tail pointer in the enqueue memory is updated it from the enqueue memory. The head pointer in the dequeue memory has to be updated it from the dequeue memory. The leaf memory is updated it with BID pointers.
8. The next FID in the FID leaf memory is compared to the tail FID in the root memory and the above operation is repeated until the next FID and the tail FID are equal. Then it is the end of the multicast group.
9. The PFQ_MEM_FULL signal is asserted to the memory manager.

Dequeue Multicast FID Algorithm

The dequeue command for a multicast FID is similar to the Unicast FID. They BID link is read from the BID memory and it will be the new head pointer. In addition the counter is read from the BID location that is the new head pointer. The counter is decremented and then written back in the same location with that BID. If the counter value is zero, then the "REL" bit in the command to the memory manager is set. Else the "REL" bit is reset and there will be no release command for that BID from the memory manager. If the BID link in the BID memory is "NULL", the queue for that FID is empty. The "NULL" has to be written back in that locations with the counter value equal to the previous count but decremented by 1.

1. The FID leaf memory is read.
2. The location read is the FID that is part of the dequeue command.
3. The head BID pointer in the FID memory will point to a location in the BID memory.
4. The location contains the next BID head pointer and the count.
5. The count is decremented by one and written into memory.
6. The new BID will replace the head pointer and the value is written in the FID leaf memory.
7. The "REL" bit in the command is set when the count is zero.
8. When the pointer is "NULL", it is handled the same way as in the Unicast FID.

Enqueue Tunneled FID Algorithm

When the queue is empty, the first enqueue is for an SOP cell. If the cell is not SOP, then the cells are dropped until the next SOP. The enqueue of the tunneled FID is exactly as the normal FID. When an SOP cell is received, the tail and head pointers are updated it as usual. When the EOP cell is received, then the FID that the PFQ sends to the DBS is not the same FID that was used to enqueue the packet. Instead, the FID will be read from the Leaf memory from a memory location pointed to by the FID of the enqueue command and it is in the ROOT Tunneling field. The FID that was read as the Tunnel ROOT memory represents one or more FIDS. The DBS will not know about the

members of the ROOT FID. The FID in order to be tunneled has to be a member of that specific tunnel.

1. The leaf FID memory is read, the tail pointer will point to the last BID in the link. The root memory association is read.
2. The current BID tail that was stored in the FID memory will be used to access the BID memory. In that location the new BID is written.
3. The new BID replaces the tail pointer and the FID location is written to the leaf FID.
4. The ROOT memory corresponding location is read.
5. If the not empty flag is reset, then the LEAF FID is written as the "CURRENT FID HEAD" and as the Tail FID
6. In EOP cell cases, the ROOT FID is sent to the DBS with the cell count.

Dequeue Tunneled FID Algorithm

The dequeue command that the DBS send to the PFQ contains an FID. The FID is the ROOT FID. The root memory has to be read and the member of the tunnel group has to be found. The Root memory has the head of the FID. The head is read and the FID leaf memory is read. After reading the Leaf memory, the dequeue operation proceeds the normal way. When the cell is EOP, then the "NEXT FID" field is used to update the new FID head pointer. If the next FID queue is empty, then the next FID field in that FID entry will be used and so on.

1. The FID Root memory is read; the current FID pointer of the tunneled FID is retrieved.
2. The FID leaf memory is read and the current BID head that was stored in the FID memory location will be used to access the BID memory.
3. The new BID will replace the head pointer and the value is written in the FID leaf memory.
4. The LEAF FID and the BID are sent to the memory manager.
5. If the BID belongs to an EOP cell, then the "NEXT FID" field is read from the leaf memory.
6. The previously read "NEXT FID" field will be used to update the "CURRENT FID TUNNEL HEAD" in the root memory.
7. The new FID head leaf memory entry is read to determine if the queue is not empty.
8. If the queue of the head FID is empty, then the "NEXT FID" field of that FID is read and the process is repeated until an FID without an empty queue is reached.
9. During
10. During the process, if the tail FID is reached, then there are no more BIDS on any queue.
11. At that point the process terminates until the next dequeue command. The "NOT EMPTY" flag is reset.
12. During the above process, assert "PFQ_ DBS_FULL" to the DBS.

Enqueue Engine

See Figure 181. The enqueue engine gets its commands from the memory. The enqueue engine issues commands to the shaper. If there are no free buffers or the request

1 is to discard the FID, then the "NULL" FID is returned. If there is an enqueue request and
 2 there are no free buffers available, then the "NULL" BID is sent to the memory manager.
 3 If the RED determines that the packet has to be discarded, then the "NULL" BID is sent
 4 to the memory manager.

5 The enqueue engine will enqueue the FID, discard the BID marked for discard,
 6 enqueue the multicast FID, and enqueue the Tunneling FID. The enqueue engine has to
 7 issue the FID and the BID memory commands always within the eight cycles window.
 8 The FID of the enqueue and de-queue commands will be compared and checked if they
 9 are the same. If they are the same FID, then only the last FID memory write is valid and
 10 has the correct information. The Enqueue information will be merged with the dequeue
 11 information and written during the dequeue operation. The BID memory accesses depend
 12 on the FID results. The data coherency of the BID during multicast dequeue operations
 13 has to be maintained by the logic with respect to the BID content of that location and its
 14 count in case of multicast.

15 The Tunneling of FIDS as well as the multicasting of FIDS is based on a circular
 16 buffer. Once the last FID is reached, then the enqueueing engine will point to the first
 17 FID in the link. Tunneling uses the root FID in the enqueue as well as in the dequeue
 18 process. The Shaper has no knowledge of the FIDS that are members of the Tunneled
 19 FID. The shaper has no knowledge of the multicast root FID.

20 When the memory manager issues a command to the enqueue engine, the FID
 21 parameters are retrieved and the BID is enqueued. If the EOP is set, then a shape
 22 command is sent to the shaper. If the FID was empty, then the FID is activated.

- 23 1. The enqueue engine will read the following:
 - 24 • FID parameters from FID enqueue memory.
 - 25 • Free buffer head from the free buffer register.
 - 26 • RED counter.
- 27 2. The enqueue engine will read the following:
 - 28 • Read the BID memory location FBH to get the new FBH.
- 29 3. The corresponding statistics from statistics memory of the FID.
- 30 4. The enqueue engine will perform the checking operations of the various
 31 conditions. The counters are incremented.
- 32 5. The enqueue engine will write the following:
 - 33 • The BID memory location current BID Tail with BID as content.
 - 34 • The BID memory location new BID tail with the "NULL" as content.
 - 35 • The "FID" enqueue location with the new parameters.
 - 36 • The "FID" dequeue location with the new tail if the queue is empty.
 - 37 • The FID statistics location with the new statistics.
 - 38 • The RED Counter with the new value.
 - 39 • The Free buffer register with the new head.

40 **Dequeue Engine**

41 The Dequeue engine gets its commands from the scheduler. The Dequeue engine issues
 42 commands to the memory manager and to the enqueue engine. When there is a dequeue
 43 command, the PFQ will get all the parameters from the FID and sends a command to the
 44 memory manager for retrieval of data. Then the PFQ will update the buffers pointers
 45

See Figure 282.

1. The dequeue engine will read the following:
 - FID parameters from the Dequeue FID memory.
 - The BID memory location BID head.
2. The de-queue engine will perform the checking operations of the various conditions. The head pointers are updated. If the head is null then the queue is empty.
3. The dequeue engine will write the following:
 - The "FID" dequeue location with the new parameters.
 - The "FID" enqueue location with the new tail "NULL" if empty.

Release BID

The free buffer engine will return the BID to the free buffer list. It can be simultaneously issued with the enqueue request. The bit "REL" will determine if the memory manager will release the BID or not.

- The free buffer BID memory with the new link BID.

Free buffer manager

There is a free buffers list. There is a corresponding head and tail registers on the device. There is a link list linked in the BID memory. The link list is initialized on init command. The enqueue will use the head of the free buffers. The tail of the free buffer list is updating it when the enqueue and dequeue engines are accessing it. The buffer number "7FFFFFFF" is reserved for not valid BID. The cell should be dropped when the memory manager gets this value in case of a discard request from the memory manger or the RED algorithm. See Figure 283.

There are three parameters that are stored in registers for the free buffer list.

- The number of available buffers. It represents the raw total storage memory space.
- The threshold of no buffers.
- The number of used buffers.

The following operations are done on the counters:

- Whenever a buffer is enqueued, then the number of used buffers is incremented.
- When a buffer is released, then the number of used buffers is decremented.
- The number of used buffers can never be less than "0".
- The number of used buffers can never be more than the available buffers.

When the number of used buffers is equal or larger to the threshold of number of buffers, then:

- All the queue of that cell has to be emptied and all the buffers of that queue are returned to the free buffers list. The buffers in use counter will be decremented by the number of cells.
- Update the free buffer register and memory link in the BID memory.

Timeout manager

It is running with simultaneously with the enqueue operation. The timeout manager is used to remove orphan cells for reassembly timeout and free the storage for incoming traffic. When an FID corresponding memory location is read for an enqueue

operation, there are two TTL bits that are reset to "00" when the FID memory location is written back. The memories that are accessed are the external FID memory and the internal leaf FID memories. The internal memory will be accessed once the external memory has been completely accessed once. The algorithm is not a fair one. When the timeout manager reads a location the TTL bits are changed from "00" to "01" or changed from "01" to "10". If the FID location is read and the TTL bits are "10", then there is a reassembly timeout and the packet remnants are discarded. The Qsize has to be decremented whenever there is a discard. It should have the value "0" in the statistics memory.

The statistics memory location is read. Every 8 cycles or every time slot, one location of the enqueue well as the statistics memories is read to perform the timeout operations. If there is a reassembly timeout, then the start of timeout operation is suspended until the data is updated in both the enqueue and the statistics memories. The timeout manager will do the following:

- An active address that points to an FID memory location. The address is sequential and circular. It is updated when the access is executed.
- Every sync cycle the FID memory is accessed, a location corresponding to that address is read and examined.
- If the TTL bits value is "00", then the value will be changed to "01".
- If the TTL bits value is "01", then the value will be changed to "10".
- If the TTL bits value is "10", then it is reassembly timeout and the packet chunk is discarded on the enqueue side.
- The TTL bits value then are reset to "00" after the discard.
- The DROP timeout bit is set. All the fragments of a packet has to be dropped until the next SOP. At the next SOP enqueueing starts again.

DISCARD FID

The discard operation is running in place of the enqueue operation. Whenever there is a discard because of Reassembly timeout, the Qsize has to be decremented. It should have the value "0" in the statistics memory. There are several discard commands:

- Memory manager generates a discard command of a packet that is being enqueued.
- RED manager generates a discard command of a packet violating the algorithm.
- Free Buffer manager generates a discard command of a packet being enqueued because of lack of free buffers. When the number of used buffers exceeds a programmed backpressure threshold register, the enqueue engine will signal the segmentation the above condition. If the number of used buffers drops 10 buffers below the programmed backpressure threshold, then the PFQ will signal the above message to the segmentation engine. If number of used buffers exceeds the programmed threshold register then the enqueue engine will discard these buffers until the next EOP. The PFQ will discard the packet fragments from the queue in such a condition.
- Timeout manager generates a discard command when the packet is not completely enqueued because the later chunks were not received or got dropped. The packet is a tail of the queue of that FID.

- 1 • CLASS discard; the packet will be discarded when the queue counter is above the
- 2 programmed threshold. On every SOP the CLASS counters are checked.
- 3 • Port Blockage; the packet will be dropped from SOP to EOP until the bit is
- 4 changed in memory. On every SOP the port register is checked.

5 When the memory manager generates the discard command, then the following steps are
6 taken:

- 7 1. The FID enqueue memory is read. The previous packet tail is parsed.
- 8 2. The BID memory is read and the BID of that location is read. The content is the
- 9 head of the new packet being enqueued.
- 10 3. The Free Buffer tail is read.
- 11 4. The BID memory location of address FBT will have the content the head of the
- 12 packet that is being discarded. The buffers are on the free buffers list.
- 13 5. All the buffers are returned to the free buffer list.
- 14 6. The FID memory location is updated with the new tail pointer, which is the tail of
- 15 the FID.

16 When the free buffer manager generates the discard command, then the following steps
17 are taken:

- 18 • The cell is dropped and the BID is returned to the free buffer list. There is no
- 19 other operation required.

20 When the timeout manager generates the discard command because of reassembly
21 timeout, then the following steps are taken:

- 22 1. The FID memory is read. The previous packet tail is parsed.
- 23 2. The BID memory is read and the BID of that location is read. The content is the
- 24 head of the new packet being enqueued.
- 25 3. The Free Buffer tail is read.
- 26 4. The BID memory location of address FBT will have the content the head of the
- 27 packet that is being discarded. The buffers are on the free buffers list.
- 28 5. All the buffers are returned to the free buffer list.
- 29 6. The FID memory location is updated with the new tail pointer.

30 When the RED manager generates the discard command, then the following steps are
31 taken:

- 32 1. The cell is dropped.

34 Port Management

35 On each SOP enqueue, the enqueue engine will check the port register to
36 determine if the packet has to be dropped or not. The port number is available from the
37 statistics memory. If the port register indicates that the cell has to be dropped, then the
38 SOP cell is dropped. If the cell is not SOP, then the drop packet flag in the FID memory
39 is checked. If the flag is set then the cell is dropped, else the enqueue operation continues.
40 The port packet discard flag for that FID is set. The flag will be reset when the port
41 register indicates that the port is re enabled again.

43 CLASS buffer Management

44 Enabling the buffer management bit in the control register will enable either the
45 RED algorithm or the CLASS management. If the RED bit is reset, in addition to the
46 previously enabled bit, then the CLASS management algorithm is enabled.

1 The enqueue engine will increment the counter in the CLASS register corresponding to
2 that FID on every cell. The Dequeue engine will decrement the counter on dequeue
3 command. When enqueue an SOP cell, the counter of the corresponding class.
4 On each SOP enqueue, the class counter will be compared to the threshold counter of that
5 class. If the counter value exceeds the threshold, then the SOP cell is dropped. The cells
6 of the rest of the packet are dropped until the next SOP cell. At that point the counter
7 value is checked against the threshold and the process is repeated.

9 **Random Early Discard Management**

10 Enabling the buffer management bit in the control register will enable either the
11 RED algorithm or the CLASS management. If the RED bit is set in addition to the
12 previously enabled bit, then the RED algorithm is enabled. Another way to bypass the
13 RED is by making both thresholds equal and very high. The average queue will not get as
14 high as the threshold. The drawback is the calculations are still needed. The queue size
15 and the averages are calculated according to cell size. The Queue size is updated every
16 time there is an EOP. It is incremented by the number of cells. The queue size is
17 decremented every time there is a dequeue operation. It is decremented by one.
18 When an SOP of an FID is received, the RED algorithm decides if the packet gets
19 discarded. If the packet is discarded, then the SOP cell is dropped as well as subsequent
20 arriving cells of the packet are dropped until the EOP cell is received. The BID is
21 returned to the Free Buffer List.

22 With constant average queue size the goal is to discard packets at regular intervals
23 in order to avoid clustering of discards and global synchronization. It is not desirable to
24 have too many discarded packets close together and also it is undesirable to have a long
25 interval between discarded packets.

26 The RED algorithm calculates the queue average size when an SOP cell is
27 received and compares it to two thresholds. If the size is less than the minimum
28 threshold, then the packet is not dropped. If the queue is larger than the maximum
29 threshold then the packet is dropped. If the average size is between the minimum and the
30 maximum thresholds, then a discard probability is calculated to determine the fate of the
31 packet. There is a random number that will be used to determine if the packet has to be
32 discarded in this case.

33 The first step when a packet is received by an empty queue is to estimate the
34 number of packets that would have been transmitted during idle time. Then the average
35 queue size will be calculated using the original estimate as if the number of packets has
36 been transmitted during that interval. The average queue size is calculated for every FID.
37 There is 1 K parameter memory for the various attributes such as the minimum and the
38 maximum thresholds and the weight of the queue. There is an association between the
39 FIDS and the attributes in the parameter memory.

40 The MAX_{th} is decreased when the probability P_a of discard is increasing. When
41 the probability changes slowly the queue average size changes slowly. This will
42 discourage oscillation in the average queue size. The values of attributes MIN_{th} and
43 MAX_{th} determines the utilization of the bandwidth

44 RED algorithm needs the following fixed and calculated parameters.
45 The following variables are calculated and saved:

- 46 • "AVG": The average queue size.

- 1 • “Q-^{TIME}”: The start of queue idle time.
- 2 • “COUNT”: The number of not discarded packets that have arrived since last
- 3 discard.
- 4 • “Q”: It is the current queue size. It is incremented on enqueue with the number of
- 5 cells and decremented with one on dequeue.

6 The following parameters are fixed and programmable:

- 7 • “W_q”: It is the queue weight. An initial value of 0.002 is used. It should not be
- 8 large in order to filter transient congestions.
- 9 • “MIN_{th}”: It is the minimum threshold for queue. If the traffic is bursty, then the
- 10 value has to be fairly large.
- 11 • “MAX_{th}”: It is the maximum threshold for queue. The value determines the
- 12 maximum average delay in the device. An initial value of 3 * MIN_{th} is used.
- 13 • “MAX_p”: It is the maximum value for a probability “P_b” of discard. An initial
- 14 value of 1/50 is used. It should never be greater than 0.1.
- 15 • “S”: Typical transmission time for a small packet.
- 16 • “RANDOM”: Random number from generator.

17 The following variables are calculated but not saved:

- 18 • “P_a”: It is the current packet discard probability.
- 19 • “TIME”: It is the current time.
- 20 • “F(t)”: It is a linear function of time t.
- 21 • “M”: The number of packets that would have been transmitted when queue is
- 22 empty.

23 The following calculations are needed:

- 24 • The number of packets “m” that might have been transmitted during idle time.
- 25 • The average size of the queue.
- 26 • The probability of discard when the queue size is between the minimum and the
- 27 maximum thresholds.

29 Equations

30 The queue empty size “M” equation is the following:

$$31 M = (TIME - Q-^{TIME})/S$$

32 The packet discard probability P_b equation is the following:

$$33 P_b = MAX_p (AVG - MIN_{th}) / (MAX_{th} - MIN_{th})$$

34 It is desirable to have P_b a negative power of 2.

35 The final packet discard probability P_a equation is the following:

$$36 P_a = P_b / (1 - COUNT * P_b)$$

37 The average queue size is calculated according to the following:

$$38 AVG = (1 - W_q) * AVG + W_q * Q$$

40 Implementation

41 It is an easy determination when the average queue size is larger than MAX_{th} and
 42 is smaller than MIN_{th}. The most complex part is when the average queue size falls
 43 between the MIN_{th} and the MAX_{th}. When the queue average size is less than MIN_{th}, then
 44 the packet is transmitted. On the other hand when the queue average size is more than
 45 MAX_{th} then the packet is discarded. If the queue average falls between the MIN_{th} and

1 MAX_{th} then the probability determines the packet destiny. A pseudo random number
2 "R" is generated by a uniform distribution random numbers generated.

3 The packet is discarded if "R" satisfies the following equation is satisfied:

$$4 R < P_b / (1 - \text{COUNT} * P_b)$$

5 If the average queue remains constant then the random number "R" is selected
6 from the uniform distribution. The N_{th} packet is discarded if the following equation is
7 satisfied:

$$8 R \leq P_b N$$

10 Random Selection

11 The random number is generated using linear feedback shift registers. The
12 Ethernet polynomial will be used to generate a random number. The number is generated
13 every time there is an enqueue operation. The following polynomial is used to generate
14 the random number.

$$15 G(X) = X^{32} + X^{26} + X^{23} + X^{22} + X^{16} + X^{12} + X^{11} + X^{10} + X^8 + X^7 + X^5 + X^4 + X^2 + X^1 + 1;$$

16 The value of the shift register is used as the number generated. The number of
17 nibbles that is stored in the attribute table is used to determine the random number is
18 used. If the value stored is "FF", then all eight nibbles are used. If the value is "7F", then
19 the first 7 nibbles are used. See Figure 284.

21 Memory Size & Total Number of BIDS

22 See Figure 285.

24 BID Memories Chip Select

25 See Figure 286.

27 EXTRA MEMORY CYCLES

28 The extra needed cycles of the FID memory are the following:

- 29 • If the queue becomes empty after a dequeue operation the tail in the enqueue
- 30 memory has to be written as "NULL".
- 31 • If the queue is empty before an enqueue operation the head in the DEQUEUE
- 32 memory has to be written with the same value as the tail pointer.
- 33 • The Dequeue operation has to decrement the Qsize in the statistics memory.

34 The BID memory requires an extra cycle after a tail is linked in memory. The extra
35 cycle will link the tail to a "NULL". When the head reaches a "NULL", then the queue is
36 empty.

38 Memories Byte Enables

39 The state machines have to drive the byte enables to the various memories. There
40 are some write cycles during which only selected byte enables have to be asserted. The
41 rest of the byte enables should not be asserted.

42 The assertion of selected byte enables during write operations eliminates the need for
43 a read operation.

- 44 • Read Operations for all memories: The state machines will assert all the byte
- 45 enables during all read operations.
- 46 • Write Operations:

- Enqueue Memory: When the queue becomes empty after a dequeue operation, the tail pointer in the enqueue memory for that FID has to be written as NULL. When writing that location, the byte enable to bits 27 downto 0 will only be asserted.
- Dequeue Memory: When the queue is empty and there is an enqueue operation, the head pointer in the de-queue memory for that FID has to be updated. When writing that location, the byte enable to bits 36 downto 0 will only be asserted because the TYPE comes with the FID.
- BID Memory: When the NULL pointer is written at the end of an enqueue operation, the BID memory location has no relevant data. All byte enables will be asserted.
- Statistics Memory: The dequeue engine has to decrement the Qsize of that FID after a dequeue operation. The STAT memory is read to get the value. Therefore, the all byte enables are asserted.

Global Synch and internal Time slot generation

The global synch is provided for all the blocks. The block has to generate locally eight time slots. The memories controllers access these memories only during certain time slots. It is important that the all tasks are completed in one cycle of eight time slots or clocks. The tasks can take longer than eight cycles to complete. But the caveat is that a task has to start every eight cycles or time slots. The latency is longer but the bandwidth is the same. The scenario complicates the data dependency of all the concurrent tasks and sequential tasks.

External and Internal Memory Initialization

On power up the memories are not initialized. There is an initialize signal that when is set by the PCI, it triggers state machine that will write a zero pattern in all memory locations of all memories except the BID memory portion of the free buffer lists. In the free buffer lists, the locations will be initialized to a value that reflects the previous BIDS. The first location of the free buffer list will have the value X"00000000". The second location will have the value X"00000001". The rest of the locations will have an incrementing values. These will be the free buffer values that are used for BIDS. The first BID "000000" is not used. It is the "NULL" pointer.

Data Dependencies and Coherency

The data and the FID of both the enqueue and dequeue previous operation that has just finished have to be saved. The current FID has to be checked against the previous ones. If the current FID is the same as the previous one of either the enqueue and dequeue operations, then the data that was saved will be used. The data in memory is not the correct one.

There are several scenarios where the data coherency will get compromised. The data coherency will get compromised with consecutive enqueue operations to same FID:

- When the consecutive enqueue operations are for the same FID the statistics counters should be updated from internal registers.

- 1 • When the consecutive enqueue operations are for the same FID the tail pointer
- 2 should be updated from internal register.
- 3 • When the consecutive enqueue operations are for the same FID the RED
- 4 Calculations should be based on data in internal registers.
- 5 • When the consecutive enqueue operations are for the same FID the Q Size should
- 6 be incremented based on data in internal register, then written in the statistics
- 7 memory.

8 The data coherency will get compromised with consecutive dequeue operations to same
9 FID:

- 10 • When the consecutive dequeue operations are for the same FID the head pointer
- 11 should be updated from internal register.
- 12 • When the consecutive dequeue operations are for the same FID the Q Size should
- 13 be decremented based on data in internal register then written in the statistics
- 14 memory.

15 The data coherency will get compromised with enqueue and dequeue operations to the
16 same FID:

- 17 • The only time that is compromised when there is an enqueue and a dequeue
- 18 operation to the same FID. If the dequeue operation renders the queue empty and
- 19 a "NULL" has to be written as the tail in enqueue memory and the head in the
- 20 dequeue memory. Instead with the enqueue operation, the queue will not be
- 21 empty any longer. The enqueue operation will preempt the dequeue operation and
- 22 the new BID will be the tail and the head pointer.
- 23 • The Qsize in the statistics memory has to be update correctly. The number of cells
- 24 is decremented by one before it is added to the Qsize.

25 The data coherency will get compromised with enqueue and timeout operations to the
26 same FID:

- 27 • The data written to the FID memory has to be the enqueue one. The TTL
- 28 increment and the discard operation have to be bypassed. There will be no
- 29 discards of the queue. The TTL will be written as "00".

30 See Figure 287.

31 **Statistics**

32 The statistics memory is accessed when there is an enqueue, dequeue and CPU
33 operations. All the requests are issued simultaneously.
34 The enqueue will update the various counters depending on the command issued by the
35 memory manager.

- 36 • Enqueue command will increment by the number of cells the Q Size counter on
- 37 EOP.
- 38 • Enqueue command will increment the packet counter on EOP.
- 39 • Enqueue command will increment the byte size counter on command.
- 40 • Discard command will increment the others discard counter.

41 The dequeue command will decrement by one the Q Size counter always.

42 The RED will increment the RED discard counter if the cell is dropped on EOP.

43

NULL Pointer

The value of "7FFFFFFF" is the "Null" pointer when the BID memories are all populated.

Shaper and Scheduler Commands flow

When it receives an EOP of an FID, then the PFQ will send a shape command to the shaper with the cell count of that packet. The Cell or packet indication has to have the correct value.

In case of ATM cells, the PFQ will send to the shaper a command to shape the FID cell.

External Interface

The external interface will be a standardized interface complying with the global SYNC. There are four interfaces between the PFQ and the external world. External memories are part of the PFQ and the interface is considered internally for all practical purposes.

- Enqueue a FID. The memory manger drives the request. The PFQ will respond with a BID and a valid signal.
 - The memory manager to PFQ command is issued in CLK 6.
 - The PFQ to memory manger response is issued in CLK 6, next synch.
- Discard a FID.
 - The memory manager to PFQ command is issued in CLK 6.
- Release BID.
 - The memory manager to PFQ command is issued in CLK 6.
- The second interface between the PFQ and memory manager is to dequeue an FID. The PFQ drives the dequeue request to the memory manager.
 - The PFQ to memory manager command is issued in CLK 6.
- There is an interface between the PFQ and the shaper. The PFQ drives the request to the shaper.
 - The PFQ to shaper commands are issued in CLK 6.
- There is an interface between the PFQ and the scheduler. The scheduler drives the requests to the PFQ. The request is to dequeue a BID of an FID.
 - The scheduler to PFQ command is issued in CLK 7.

Memory Structure**RED Attribute Memory**

This is internal memory. There are 1K entries of this memory. The software writes it only. The address is the "RED_ASSOCIATION" from the FID memory. See Figures 288 and 289.

Statistics Memory

This is an external memory. There are 2 MEG locations. There are 2 locations per FID. The software writes only the port number. The rest is written by the logic. The CPU reads the counter values. The address is the "FID" for bits 20 downto 1 with bit 0 as "0" or "1".

See Figures 290, 291 and 292..

All discard counters are updated on EOP. If the packet is being dropped because of RED and the EOP cell has a discard command instead of enqueue, then the RED counter is not update it.

FID Dequeue Memory

This is an external memory. There are 1 Meg locations. There is one location per FID. The software does not write any location. The software under normal operating traffic conditions does not need any values from this memory. The address is the "FID". See Figures 293 and 294.

FID Enqueue Memory

This is an external memory. There are 2 MEG locations. There are 2 locations per FID. The software writes only the RED association. The software under normal operating traffic conditions does not need any values from this memory. The address is the "FID" for bits 20 downto 1 with bit 0 as "0" or "1". See Figures 295, 296 and 297.

Buffer ID Data Structure

This is an external memory. There are 8 MEG locations. There is one location per BID. The software writes the count only. The count cannot change when the queue is not empty. Hence there is no need for a separate byte write enable. The address is the "BID" tail or head that are read from the dequeue memory, enqueue memory or the free buffer register. See Figures 298 and 299.

The SOP, EFCI, OAM are set when the SOP cell arrives. The EOP is set in a different BID when it arrives. There is no need to write extra cycles on EOP.

Multicasting Data Structure

These are internal memories.

Multicasting FID LEAF Data Structure

Multicasting FID Dequeue Memory

There are 1 K locations. There is one location per FID. See Figures 300 and 301.

Multicasting FID Enqueue Memory

There are 2 K locations. There are two locations per FID. There are 2 locations per FID. The software writes only the RED, and "NEXT_TUN". The software under normal operating traffic conditions does not need any values from this memory. See Figures 302, 303 and 304.

1 **Multicasting FID ROOT Data Structure**

2 The enqueue engine uses the next table. There are 1 K locations. The software
3 writes the head and the tail. See Figure 305.

5 **Tunneling Data Structure**

6 The memories are internal.

7

8

9

10

11 **Tunneling FID LEAF Data Structure**

12 ***Tunneling FID Dequeue Memory***

13 There are 1 K locations. There is one location per FID. There is one location per
14 FID. The software does not write any location except the "TUN_ROOT". The software
15 under normal operating traffic conditions does not need any values from this memory.
16 See Figure 306.

17

18 ***Tunneling FID Enqueue Memory***

19 There are 2 K locations. There are 2 locations per FID. There are 2 locations per
20 FID. The software writes only the RED, and "NEXT_TUN". The software under normal
21 operating traffic conditions does not need any values from this memory. See Figures
22 307, 308 and 309.

23

24 **Tunneling FID ROOT Data Structure**

25 The de-queue engine uses the next table. There are 1 K locations. The software
26 writes the head and the tail. See Figure 310.

27

28 **Interface**

29 The interface is either internally to the memory manager, shaper or the scheduler.
30 Else, it is to the external rams. The BYTE enable for the enqueue memory should be
31 connected that the "BYTE_ENABLE_0" is connected to the first three bytes that contain
32 the BID tail pointer. The second "BYTE_ENABLE_0" should be connected to the rest of
33 bytes. When writing the BID tail pointer with the value "NULL" by the dequeue engine,
34 there is no need for a read operation, instead only the first three bytes are written. The
35 rest of the data field is maintained as it was. The dequeue memory has no relevant data
36 that need to be maintained when the enqueue engine updates the BID head pointer after
37 empty queue. The type field is updated with the pointer. The "TYPE" is part of the

1 command. See Figure 311. The signal "MEM_PFQ_FULL" is synchronized and bussed
2 through to the DBS.

3 4 **CPU Interface**

5 There are signals that are generated by the CPU interface and are used to control
6 the memory sizes and the enable and disable signals. See Figure 312.

7 8 **Segmentation Interface**

9 See Figure 313. The PFQ of an egress device sends a message to the segmentation
10 engine of Ingress devices when the number of used buffers exceeds a programmed
11 threshold. The data will contain the port address of the egress device. The segmentation
12 engine will multicast the message to all the input devices that instructs these devices to
13 refrain from sending any traffic to the device of the corresponding port. See Figure 314.
14 The message indicates the condition when the number of buffers used exceeds the
15 threshold value by setting the flag to '1'. There is another message that indicates when
16 the condition no longer exists by setting the flag to "0". The rest of the data bits indicates
17 the port number that is programmed by the CPU.

18 19 **DBS & CBWFQ Interface**

20 See Figure 315.

21 **DBS Interface**

22 See Figure 316.

23 24 **Memory Manager Interface**

25 A command is stretched to few cycles as described below so that the number of
26 pins is reduced. The interface between the 2 blocks is shown in Figures 317 and 318.
27 The En-queue/Discard command from the memory manager to the PFQ takes 2 cycles to
28 complete. The De-queue command from the PFQ to the memory manager takes 3 cycles
29 to complete. When the BID is not valid, the "NULL" BID will be returned to the memory
30 manager as the BID value. The memory manager will drop the cell whenever the BID is
31 "NULL". See Figures 319 and 320.

32 33 **MEM to PFQ**

34 The memory manager starts the command by asserting valid signal. The PFQ will
35 always sample the command. If it is not NOP, then the command determines the number
36 of valid cycles with data. If the command is a "discard" command, then the only valid
37 parameters are the FID; all the rest of the parameters should be ignored by the PFQ. If the
38 command is a release then the BID is important and the rest is ignored. If it is an enqueue
39 command, then both D1 and D2 are valid. See Figure 321.

40 The memory manager as well as the PFQ should have FIFOS to store the
41 commands and the parameters.

42 43 **PFQ to MEM**

44 The command takes 3 cycles to complete. In this command the data fields are
45 transferred in 3 cycles as shown in Figure 322.

If "REL" is set, then the memory manger can reuse the buffer, which is the BID. As a result, the memory manger cannot send a command to PFQ releasing BID that was dequeued with the "REL" bit reset. The PFQ will not check the "DEQ CNT" in the BID memory. The "DEQ CNT" is used for multicasting traffic multiple times.

Interface Timing between the Memory Manger and PFQ

The MX1 and MX2 are not synchronous. The global synch signal is generated locally. As a result, the time slot is not adhered to in this interface. All the commands or data are stored in FIFOS after they reach the receiving device. The response to such a request will take between 16 and 24 clock cycles. For instance, the memory manger needs to wait between 16 and 24 clocks before it gets a BID after it asserts the FID enqueue command. See Figure 323.

Memory Manager Interface Timing

See Figure 324. The clock is 100 MHZ and the commands from the PFQ to the memory manager are valid every 4 clocks. Similarly, the commands from the memory manager to the PFQ are valid every 4 clocks. As a result, the PFQ block and the memory manager block will receive a command every 8 clocks of 200 MHZ.

SEGMENTATION Interface Timing

See Figure 325. The clock is 100 MHZ and the data is serially driven on two data pins. The flag indicates if the queue is full when it is "1". When the queue is full and the number of buffers used is equal or exceeds a programmable threshold, the PFQ will send a message indicating the condition by setting the flag to "1". When the number of buffers used is less than the programmed threshold by 10, then an equivalent message is sent but with the flag reset.

Software Interface

The CPU interface doesn't use bursting. Data bus is still 32 bits.

Features

- Data bus of 32 bits. There is data bus in and out.
- Single address access (not a burst access).
- CPU accesses all registers in all blocks via direct read/write commands.
- All external memories and most of the internal memories are accessed via indirect accesses through registers
- CPU accesses all internal/external memories in all blocks via indirect accesses using the general-purpose registers dedicated for these accesses.

Address space

The address space is 64 entries only. The internal registers directly mapped to these addresses and external/internal memories are indirectly mapped. The address spaces per lookup engine is shown in Figure 326.

Direct/Indirect Accesses

A direct access is used to access. If there are internal memories or registers, which are longer than 32 bits, then read or write to these locations will be executed in more than one access. For example an entry in internal memory of 50 bits width (or a register with the same length) will be accessed in two CPU commands first access bits [31:0] and then [49:32]. Each register or small internal memories have their own address, which is unique to the entire chip set.

An indirect access is used to access internal and external memories. Each block has a set of registers dedicated for the indirect access. One register contains the address and the command (Rd, Wr, Init etc). Few other registers contain the data for a write command. Writing to the COM/ADD register triggers the execution of the command, once the block completed the execution of the command, the COM portion of the command is cleared by the block. When a read command occurs the data is written to the registers and then the CPU reads the registers to get the data.

Registers

See Figure 327.

Commands**NOP.**

See Figure 328.

Read FID Enqueue Memory (72 bits) – Up to 2M locations.

See Figure 329.

Write FID Enqueue Memory (72 bits) – Up to 2M locations.

See Figure 330.

Read FID Dequeue Memory (36 bits)– Up to 1M locations.

See Figure 331.

Write FID Dequeue Memory (36 bits) – Up to 1M locations.

See Figure 332.

Read BID Memory (36 bits) – Up to 16M locations.

See Figure 333.

1 **Write BID Memory (36 bits) – Up to 16M locations.**

2 See Figure 334.

3 **Read Stat Memory (72 bits) – Up to 2M locations.**

4 See Figure 335.

5 **Write Stat Memory (72 bits) – Up to 2M locations.**

6 See Figure 336.

7 **Setup connection command for FID**

8 See Figure 337. This command sets up a connection by accessing external
9 memories and initializes the corresponding locations to the proper parameters. It is a
10 write command. Clear PFQ entries for ENQ, DEQ and STAT!

11

12 **Tear-Down connection command for FID**

13 See Figure 338. This command tears down a connection by accessing external
14 memories. It is a write command. The data is all "0"

15 PFQ returns all remaining BIDS on the queue of the down connection to the free
16 link list during this command.

17

18 **INIT FID MEMORIES**

19 The memories are the enqueue, dequeue, and statistics memories. See Figure 339.

20

21 **INIT BID Memory**

22 This will build the free buffer list. See Figure 340.

23

24 **Reserved Op-codes**

25 The op-codes 1101, 1110, and 1111 are reserved.

General commands

The following commands are combinations of the above simple commands. A user command like this should be broken to several Read or Write commands in the driver.

Initialized External memories

- Write "0" in all memory locations.
- Write "NULL" in the enqueue memory.
- Write "NULL" in the dequeue memory.
- Build the free buffer list by linking each buffer to the next one.

Initialized Internal memories

- Write "NULL" in the Multicast Root memory.
- Write "NULL" in the Multicast Leaf memory.
- Write "NULL" in the Tunneling Root memory.
- Write "NULL" in the Tunneling Leaf memory.

Principle Of Operation

The PFQ performs the following operations:

- PFQ adds a BID to the FIDS queues. In essence it enqueues a BID.
- PFQ removes a BID from the FIDS queues. In essence it dequeues a BID.
- PFQ adds a BID to the free buffer queue. In essence it de allocates or returns a free BID.
- PFQ removes a BID from the free buffer queue. In essence it dequeues or allocates a free BID.
- PFQ increments the statistics counters.
- PFQ timeouts the active FIDS.
- PFQ asserts commands to the shaper and memory manager.
- PFQ accepts commands from the scheduler and memory manager.
- PFQ transfers statistics to the CPU.
- PFQ executes early random discards depending on criteria that associate the FID with a counter and an algorithm.
- PFQ executes discard operations.
- The Memory manager issues the enqueue, discard and the release.
- The scheduler issues the dequeue command.

The Per flow queue will get a command from the memory manager to enqueue an FID. The PFQ rejects or honors the request depending on the free buffers available or RED. The request is honored and the enqueue operation is performed. If the request is rejected, then a discard operation is performed by making the tail of the discarded packet

as the tail of the free buffer list. Also, the head of the packet will be linked to the previous free buffer tail.

The scheduler commands are dequeue ones. The PFQ will remove the BID from the queue and transfer the information to the memory manager.

Error Conditions

The total numbers of free buffers is compared to the number of buffers used and also to the total number of buffers on the queues. If the numbers do not match, then there will be an abort of all the operations. The PFQ will indicate the condition to the CPU. Another error condition happens when there is a dequeue request and the queue is empty.

TIMING DIAGRAMS

The timing diagrams are for FID and BID memory accesses. They include the accesses of the FID memory as well as the BID memory for enqueueing, dequeueing, statistics collections and timeout scanning. See Figure 341.

The memories cycles are summarized in Figures 342-346.

When dequeues engine is updating the tail pointer in the enqueue memory to "NULL", there is one write operation that is performed to the enqueue memory location 0. The BYTE ENABLE for the first 24 bits should be used. All the other byte enables should not be asserted.

When the enqueue engine is updating the head pointer in the enqueue memory to "BID", there is one write operation that is performed to the dequeue memory. The BYTE ENABLE for the first 24 bits should be used. All the other byte enables should not be asserted.

Per Flow Commands

- The Memory manager command to PFQ:
 - ENQ {FID}: Enqueue the FLOW.
 - DSCD {FID}: Discard the FID.
 - REL {BID}: It releases the BID.
- The PFQ Commands to the memory manager:
 - SND {BID}{REL}{PORT}: Get the data cell belonging to BID and send it with port number "PORT", afterwards release the buffer if "REL" is set.
 - VALID {BID}: Accept the BID for the previous enqueue command.
- The PFQ Commands to DBS:
 - SHAPE_OR_SCHEDULE {FID}{BID}{LOCATION}{#CELLS}{EOP}: Shape the current Flow ID FID with the location "LOCATION", EOP marker received even though should be implied.
- The DBS Commands to PFQ:
 - DEQ {FID}: The PFQ will pass the Buffer ID to the memory manager. If EOP is set, then send a shape command to shaper.
- The PCI Commands to PFQ:
 - Get [ADR]{CNTS}: Get the counts values from the statistics memory.

Data Base

The data base block is responsible for managing the packets link list per FlowID for FlowIDs in the Shaper and FlowIDs in the scheduler. The block receives an En-queue command from the PFQ and sends the FlowID either to the Shaper or to the Scheduler. A De-queue command starts with a port selection based on a port calendar and then a FlowID is selected between the scheduler and the shaper.

Shaper traffic has priority over scheduler traffic. Shaper traffic is divided to 8 classes of strict priority, once all shaper's classes are empty for a given port, the scheduler is allowed to schedule.

Features

1. This block operates at 200MHz clock.
2. Support up to 1M FlowIDs.
3. Support a link list of packets/cells up to 8M cells.
4. Performs the following in one port slot (8 cycles):
 - a. One input phase from the PFQ and redirect to an input phase in the Shaper or in the Scheduler or to both (when traffic is just metered)
 - b. One output phase of the Shaper.
 - c. One output phase of the Scheduler.
5. Support 8 strict priority classes for the shaped traffic per port.
6. Performs port scheduling for the shaped traffic base of strict priority within the block.

Interface

A global block diagram of the database block is shown in Figure 347. Figure 348 shows the interface to/from the data base block. See also Figure 349.

Functional Description

The DBS block acts as the holder of the information per FlowID and the packets within the FlowID for the Shaper and the Scheduler blocks. The DBS block can be divided to several processes during one port slot:

PFQ Input Phase

For a packet traffic the PFQ first assembles a complete packet and then transfer the packet to the DBS block, when in a cells traffic, a cell is being sent to the DBS each time it arrives. One cell in cells traffic is considered as a packet with a length of one cell. When a packet arrives to the DBS block from the PFQ, the DBS first links the packet to the FlowID and also generates an input phase to the Shaper or to the Scheduler or to both based on a parameter bit per FlowID.

Shaper Output Phase

The output phase of the Shaper is being performed on a cell basis, means each cell of a packet generates an output phase from the Shaper. Once a cell is driven from the Shaper, the DBS block decrement the cell counter for the specific packet. If it's the last cell of the packet, the DBS links the packet to the shaper's output memory per port. The

DBS also generates an empty indication to the Shaper if there are no more packets to shape for the FlowID.

Port Calendar

A port calendar is a free running port selection based on a memory programmed by the user. Once a port is selected for the current output phase, the port number is transferred to both scheduler and shaper's output memory. Both blocks select their output FlowID for the output stage. If the shaper's output memory has a valid FlowID and the scheduler is not in the middle of a packet, then the FlowID to be driven to the PFQ is always from the shaper's output memory. Shaped traffic has priority over un-shaped traffic.

Each output phase that starts with the port calendar also starts with checking the port status. Only if the port calendar has a valid port and that port is not full according to the port status, a valid port is driven to the scheduler and to the shaper.

Port Status

This block receives the status information from the output SPI interface for 64 ports plus one CPU port. See the timing diagram of that interface below.

Scheduler Output Phase

Once a port is selected, the scheduler performs it's own output phase. The DBS block needs to support this phase the same way it supports the shaper's output phase.

Block Description

The following are the algorithms to perform in each of the previous phases:

PFQ Output Phase

1. The PFQ drives the following signals: *VALID, FID, NCELL, BID, CLPI*

2. *Read FID_mem[FID] { rp, wp, e, cell_cnt, ncell, clp, port, qos, c_p, CPU_port, shape}*

3. *n_wp <= BID*

4. *n_e <= 0*

5. *if(e==1)*

1) *if(shape_meter==0 | shape==1) perform input phase to Shaper:*
VALID, FID, qos, ncell

2) *if(shape_meter==0 | shape==0) Perform input phase to Scheduler:*
VALID, FID, qos, CPU_port

3) *n_rp <= BID*

4) *n_ncell <= NCELL*

5) *n_cell_cnt <= NCELL*

6) *n_clp <= CLPI*

7) *pkt_next_we = 0*

6. *else*

1) *if(shape_meter==0) perform input phase to Shaper:*
VALID, FID, qos, ncell

2) *n_rp <= rp*

- 3) $n_ncell \leq ncell$
- 4) $n_cell_cnt \leq cell_cnt$
- 5) $n_clp \leq clp$
- 6) $pkt_next_we \leq 1$
7. if($VALID==1$) Write $FID_mem[FID] \leq$
 $\{n_rp, n_wp, n_e, n_ncell, n_cell_cnt, n_clp, qos, c_p, CPU_port, shape\}$
8. if($VALID==1$ & $pkt_next_we==1$) Write $pkt_mem[wp] \leq$
 $\{ncell, clp, FID, next, [dbs_dpkt_word_en[1:0] = 11] \leq \{NCELL, CLPI,$
 $0, BID\}$

Shaper Output Phase

1. The SHP drives the following signals: $VALID, FID, VIOLATOR$
2. Read $FID_mem[FID]$ $\{rp, wp, e, cell_cnt, n_cell, clp, port, qos, c_p,$
 $CPU_port, shape\}$
3. Wait until read is complete
4. $cur_eop \leq cell_cnt == 1$
5. $cur_empty \leq cur_eop \& rp == wp$
6. Read $pkt_mem[rp]$ $\{ncell, clp, next, FID\}$
7. Read $shp_out_port[port]$ $\{p_rp, p_wp, p_e\}$
8. if($cur_eop == 1$) $n_rp \leq next; n_cell_cnt \leq ncell; n_n_cell \leq ncell;$
9. else $n_rp \leq rp; n_cell_cnt \leq cell_cnt - 1; n_n_cell \leq n_cell;$
10. if(cur_empty) $n_e \leq 1$
11. $n_FID \leq FID$
12. $n_clp \leq en_clp_change ? (VIOLATOR | clp) : clp$
13. if(cur_eop)
 $shp_out_port_we \leq 1$
 $n_p_e \leq 0$
 $n_p_wp \leq rp$
 $if(p_e) n_p_rp \leq rp$
14. else
 $shp_out_port_we \leq 0$
15. send cur_empty indication to Shaper
16. if($VALID==1$) Write $FID_mem[FID] \leq$
 $\{n_rp, wp, n_e, n_cell_cnt, n_n_cell, qos, c_p, CPU_port, shape\}$
17. if($VALID==1$ & cur_eop)
 $Write\ shp_out_port[port] \leq \{n_p_rp, n_p_wp, n_p_e\}$
18. if($VALID==1$ & $cur_eop \& p_e == 0$)
 $Write\ pkt_mem[p_wp][next][dbs_dpkt_work_en[1] = 1] \leq \{rp\}$
19. if($VALID==1$ & cur_eop)
 $Write\ pkt_mem[rp][dbs_dpkt_work_en[0] = 1] \leq \{n_clp, n_cell, n_FID\}$

DBS Output Phase (Combined with Scheduler Output Phase)

1. Read $port_calendar\ mem[count]$ $\{portid, valid, jump\}$
2. if($jump==1$) $count \leq 0$
3. else $count = count + 1$
4. send shaper and scheduler: $portid, valid$

```

1      5. wait for 4 cycles to receive SHP_FID, SHP_VALID, SCH_FID, SCH_VALID,
2      SCH_MID
3      6. get PFQ_FULL indication from PFQ block
4      7. if(valid & ~PFQ_FULL & SHP_VALID & (~SCH_MID | ~SCH_VALID))
5          1) out_from_shp <= 1
6          2) out_from_sch <= 0
7      if(valid & ~PFQ_FULL & SCH_VALID & (~SHP_VALID | SCH_MID))
8          1) out_from_shp <= 0
9          2) out_from_sch <= 1
10     8. if(out_from_shp==1)
11         1) Read shp_priority_mem[portid] {empty}
12         2) Select class per portid based on empty : class_addr = {portid, class}
13         3) Read shp_out_mem[class_addr] {p_rp, p_wr, p_e}
14         4) Read pkt_mem[p_rp] {ncell, clp, next, FID}
15         5) Send to PFQ: FID, portid, p_clp, valid(=1)
16         6) If(ncell==1)
17             a. n_p_rp <= next
18             b. if(p_rp==p_wp) n_p_e <= 1
19             c. n_ncell <= 0
20         7) else
21             a. n_p_rp <= p_rp
22             b. n_p_e <= p_e
23             c. n_ncell <= ncell - 1
24         8) Write shp_out_mem[class_addr] <= {n_p_rp, p_wp, n_p_e}
25         9) If(n_p_e==1) Write shp_priority[portid] <= new empty indication
26         10) Write pkt_mem[p_rp] <= {n_ncell}
27     9. if(out_from_sch==1)
28         1) Read FID_mem[SCH_FID]
29     {rp, wp, e, cell_cnt, n_cell, port, qos, c_p, CPU_port, shape}
30         2) eop <= (cell_cnt==1)
31         3) empty <= (cell_cnt==1 & rp==wp)
32         4) busy <= PFQ_FULL
33         5) Read pkt_mem[rp]{ncell, clp, next, FID}
34         6) Send eop, busy, empty to scheduler
35         7) if(eop) n_rp <= next; n_cell_cnt <= ncell; n_n_cell <= ncell;
36         8) else n_cell_cnt <= cell_cnt-1; n_n_cell <= n_cell; n_rp <= rp;
37         9) if(empty) n_e <= 1
38         10) Write FID_mem[SCH_FID] <=
39     {n_rp, wp, n_e, n_cell_cnt, n_cell, port, qos, c_p, CPU_port, shape}

```

Port Status

This interface drives the status of the 64 ports in the output SPI and the CPU port. It is a continuous interface. Every 8 cycles of 200MHz the status indication of the 65 ports is updated to a 65 bits register. The first cycle of the eight is marked with a specific signal (sync signal). Figure 350 shows the interface.

The database block needs to synchronize the incoming status bits to the local 200MHz clock and then use the status bits. A synchronization FIFO is being used, but since the write and read frequencies are the same, there is no empty or full indications. Reads will always lag 3 cycles after the writes. To achieve that, the read pointer is always the value of the write pointer minus 3. Also the first 3 reads after reset should be ignored.

Full indication from PFQ

A full indication from the PFQ can be asserted any time. Once a full indication is asserted, the DBS first complete the current output phase and then freeze the increments of the port calendar. From that moment until this indication is de-asserted no output phase is performed from the scheduler or from the shaper's output FIFOs. The shaper's output phase to the DBS should still continue without any connection to the full indication.

CLP Bit

The CLP bit is modified in the Database and then is sent back to the PFQ during the output stage. Only FlowIDs that are shaped can have their CLP bit changed, FlowIDs to the scheduler have the CLP bit unchanged. The violation indication from the shaper should set the CLP bit to 1, if no violation indication arrives then the CLP bit is to remain the same as it was.

Memory Description

Chip-Select indication for the PKT Memory

The packet memory can support up to 8M locations (the address bus has 23 bits). Since there are no SSRAM memories in the market that have that amount of entries, the DBS provides chip-select signals so that several memories can be connected to get up to 8M locations.

The chip-select signals are used only when supporting 2M locations and up. The generation of the chip-select depends on the input from the CPU interface called CPU_CS_SEL[1:0] as follow:

- 00 – generate chip-select based on address[20:19]
- 01 – generate chip-select based on address[21:20]
- 10 – generate chip-select based on address[22:21]
- 11 – Reserved

for example, if supporting 2M locations only and the memory devices used have a depth of 512K locations, then CPU_CS_SEL bus should have the value of 00. the chip_select[3:0] will have the following meaning:

- chip_select_l[0] – selects memory block for addresses 0 – 1/2M
- chip_select_l[1] – selects memory block for addresses 1/2M – 1M
- chip_select_l[2] – selects memory block for addresses 1M – 1 1/2M
- chip_select_l[3] – selects memory block for addresses 1 1/2M – 2M

In this case address bus bits [18:0] are going to all memory blocks (Note, each block has only 512K location, which translated to 19 bits), the chip_select_l bus controls 4 memory blocks, therefore, the total locations supported is 2M.

1
2 **DBS FlowID Memory**

3

4 **FlowID Memory**

5 Address: FID, Length: 1M, #bits: 91

6 See Figure 351. The user has to initialize the above fields, which are marked as
7 software using setup connection command.

8

9 **PKT Memory**

10 Address: BID, length: 8M, #bits: 55

11 See Figure 352.

12

13 **Port calendar memory**

14 Address: counter, length: 192, #bits: 8

15 See Figure 353. The user should initiate this memory prior to any traffic based on the
16 users input.

17

18 **Shp_strict memory**

19 Address: port, length: 64, #bits: 12

20 See Figure 354.

21

22 **Shp_out_port memory**

23 Address: port, length: 512, #bits: 49

24 See Figure 355

25 .

26 **Registers Description**

27 See Figure 356.

28

29 **CPU Commands**

30

31 **Read external FID memory**

32 See Figure 357.

- 1 A[0] – C_P (bit [6] in the register)
- 2 A[1] – CPU_PORT (bit [7] in the register)
- 3 Pri[2:0] – Priority for shaped traffic
- 4 S – Shape (bit [23] in the register)
- 5 C[0] – clpi (bit [10] in the register)
- 6 C[1] – Empty (bit [11] in the register)

7

8 **Write external FID memory**

9 See Figure 358.

- 10 A[0] – C_P (bit [6] in the register)
- 11 A[1] – CPU_PORT (bit [6] in the register)
- 12 Pri[2:0] – Priority for shaped traffic
- 13 S – Shape (bit [23] in the register)
- 14 C[0] – clpi (bit [10] in the register)
- 15 C[1] – Empty (bit [11] in the register)

16

17 **Setup FID connection**

18 See Figure 359.

- 19 A[0] – C_P (bit [6] in the register)
- 20 A[1] – CPU_PORT (bit [6] in the register)
- 21 Pri[2:0] – Priority for shaped traffic
- 22 S – Shape (bit [23] in the register)
- 23 During the setup connection command, the user defines the above fields only.
- 24 The fields RP, WP, NCELL, CELL_CNT and CLP should be reset to 0 by the hardware
- 25 The field EMPTY should be set to 1 by the hardware.

26

27 **Read external PKT memory**

28 See Figure 360.

- 29 L[0] – CLP
- 30 L[1] - Reserved

31

32 **Write external PKT memory**

33 See Figure 361.

- 34 L[0] – CLP

1 L[1] - Reserved

2

3 **Read port calendar memory**

4 See Figure 362.

5 Data[5:0] – portid

6 Data[6] – valid

7 Data[7] – jump

8

9 **Write port calendar memory**

10 See Figure 363.

11 Data[5:0] – portid

12 Data[6] – valid

13 Data[7] – jump

14

15 **Read SHP Strict memory**

16 See Figure 364.

17 P[2:0] – previous QOS pointer

18 P[3] – previous QOS pointer valid

19

20 **Write SHP Strict memory**

21 See Figure 365.

22

23 **Read SHP out memory**

24 See Figure 366.

25

26 **Write SHP out memory**

27 See Figure 367.

1

2 **Init SHP_out_port memory**

3 See Figure 368. This command generates write accesses to all entries in the
 4 Shp_out memory (64 entries) with a specific data of: p_rp = 0, p_wp = 0, p_e = 1

5

6 **Init Shp_strict memory**

7 See Figure 369. This command generates write accesses to all entries in the
 8 Shp_strict memory (64 entries) with a specific data of: empty = 8'b1111_1111,
 9 prev_qos_v = 0, prev_qos = 0

10

11 **Reserved Opcodes**

12 The opcodes: 1110, 1111 are reserved.

13

14 **Timing Diagram**

15 Output phase to PFQ from SHP output memory. See Figure 370.

Notes for Figure 370:

PFQ Output, SHP Output, SCH Output DBS Output

Notes :

1. Slot 4 has data dependency on slot 6 if SHP FID and PFQ FID are same. Use following algorithm for this data dependency

if (shp_FID == PFQ_FID & shp_out_valid & PFQ_in_valid)

@slot 6 empty <= 0

wp <= BID(PFQ)

if (shp_empty) rp <= BID(PFQ)

else rp <= n_rp

2. Slot 4 has data dependency on slot 7 if scheduler FID is same PFQ FID. In that case

if (sch_FID == PFQ_FID && sch_out_valid & PFQ_in_valid)

@slot 7

empty <= 0;

wp <= BID(PFQ)

if (sch_empty) rp <= BID(PFQ)

else rp <= n_rp;

There are no data dependencies on the pkt memory since a packet can exist only in one place at a given time.

16

1 Interfaces

3 Shaper

4 See Figure 371.

6 Scheduler

7 See Figure 372.

9 PFQ

10 See Figure 373.

12 Output Scheduler

13 The output scheduler receives unshaped flow IDs from the Data Base Module(it is
14 referred as DBS) and links them in its database. Then it chooses a flow ID based on its
15 scheduling algorithm and outputs the flow ID to the PFQ engine. Note that unshaped
16 flow IDs will get available bandwidth which is not used by shaped FLOW IDs. Shaped
17 flow IDs are scheduled for output by the DBS and these flow IDs get committed
18 bandwidth. In other words Shaped flow IDs from the DBS always get priority over
19 unshaped flow IDs from scheduler unless it is middle of packet. When DBS received a
20 flow ID, flow ID is either sent to shaper(if it to be shaped) or scheduler(if it is not to be
21 shaped). Once flow ID is received from DBS, scheduler links this flow ID based on the
22 port number and QOS. Ports are always serviced in round robin fashion by the database.
23 DBS informs the scheduler the port number which will be served in this time slot(8
24 clocks). If there are one or more flow IDs pending for this port then a flow ID is selected
25 for scheduling based on the following criteria.

26 The scheduler implements 2 levels of scheduling

- 27 • Level 2: Weighted-Round-Robin algorithm combined with strict priorities and
28 best effort to schedule QOSs.
- 29 • Level 3: Round Robin between FlowIDs

30 The scheduler also interfaces with a CPU through the software interface. The CPU
31 can read or write the scheduler internal registers and all the memories. Furthermore, the
32 scheduler receives flow control status from the re-assembly engine, the PFQ engine and
33 the egress chip.

35 Features

- 36 1. The scheduler operates at 200MHz clock.
- 37 2. Round-robin algorithm between Flows.
- 38 3. Combination of weighted round-robin algorithm and strict priority for QOSs.
- 39 4. Support up to 1048576 Flows (1M).
- 40 5. Support 192 port of STS1 rate each up to one OC192 port. CPU port is always OC3.

6. Supports up to 1024 QOSs that are shared between the output ports.
 - Up to 8 QOSs per port. CPU port has only one QOS, which is not accounted in the 1K.
 - The 8th QOS per each port is acting as a best effort QOS, means this QOS will be selected only if all other QOSs are empty.
7. Two modes of operations:
 - Egress Mode: (Normal operation) Each port has up to 8 QOSs, there are 64 output-ports plus 1 CPU port which has only one QOS.
 - Ingress Mode: Support only one output port. All QOSs belong to that one port. Virtually there are still up to 256 ports with up to 8 QOSs dedicated to each. A separate port status information bits will arrive from the Egress chip to block the QOSs in the Ingress chip, which send traffic to the congested ports in the Egress chip. This mode of operation affects the port calendar only.

Figure 374 shows both applications with respect to the switch:

8. Status for Ingress chip applications comes from the Egress chip.
9. External memory supported is pipelined ZBT memory. Micron provides pipelined ZBT memories of 512K x 32/36 or 1M x 18.
10. In case of scheduling packets, only one packet can exist in the scheduler for a specific FlowID. Once this packet is scheduled the database block (DBS) will issue another packet of the same FlowID to the scheduler if there is one linked in the DBS. In case of scheduling ATM cells, each cell is treated as a one cell packet.
11. Accept to link a flow from the DBS once every 8 cycles of 200MHz.
12. Since 8 clock cycles are needed to service a port (a cell), a total of 520 (64 ports + CPU port) cycles are needed to keep up the 10Gbs traffic rate.

Interface

Figure 375 shows the interface to/from the scheduler block. Figure 376 shows the detailed interface to other blocks. Signal Descriptions are found in Figure 377.

Functional description

Flow ID is given to scheduler by DBS if flow ID given by PFQ is not to be shaped. Based on a port and a QOS a flow ID belongs to, the scheduler informs DBS to link up the flow ID in the DBS's linked list memory structure. DBS selects port to serve and inform scheduler port number which will be served in the given slot. Scheduler then selects a QOS based on the strict priority or WRR algorithm among all the QOSs of the port, and finally chooses a flow ID based on the round robin algorithm between all the flows of the QOS. The scheduler sends the scheduled flow ID to DBS for scheduling. At this point DBS decides whether flow ID from scheduler will be sent out or flow ID from shaper will be sent out. Always if there is flow ID from shaped traffic for this port, it will be sent out first. In that case DBS informs scheduler not to update any parameters as flow id from scheduler has not been scheduled. When there is no flow ID from shaped traffic for this port, then flow ID from scheduler will be sent out. In that case, scheduler updates all parameters. Also if the current cell scheduled is the last cell (in case of ATM traffic, Every cell sent out will be marked as eop) then End Of Packet(eop) signal is sent from DBS to scheduler. Also if the current packet is the last packet linked for this flow ID,

empty will be asserted by DBS to scheduler. If it eop and not empty then scheduler has to rotate flow IDs. Scheduler will link current read pointer(f_rp) to end of link(f_wp). This information is sent to DBS. Also next_FID given by DBS will become head of link list for this QOS.

The following explains all the scheduling algorithms employed by the scheduler and the scheduler's linked list memory structure.

Round Robin Algorithm

All the elements to be scheduled are linked together in one link list. There are two pointers to that link list, a head pointer and a tail pointer. A selection is basically forwarding the head pointer. There are 2 options to continue the algorithm after selecting the element: the element is linked to the end of the list and become the tail pointer while the next element becomes the head of the list. The second option is that the element is removed from the link list, in this case only the pointer is changed to the next element in the list. See Figure 378.

A Round Robin algorithm is used for the FlowIDs. The port calendar is implementing the same but without removing or adding elements to the link list, means a fixed list.

Weighted-Round-Robin Algorithm

This algorithm is an extension to the Round Robin one. Each element in the link list has a weight and will be selected as much as the weight allows. This algorithm should generate a relative bandwidth to each element according to the total elements linked:

$$BW_k = \text{Weight}_k / \sum \text{Weight}_i$$

To implement this algorithm there is a need for two link lists. Active link list, which holds the active elements to schedule, once an element k was schedule weight_k times, it is moved to the waiting list. Once all the elements in the active list exhausted their weight and moved to the waiting list (active list becomes empty) the waiting list becomes the active list and the above formula is performed on the elements in the new active list. During processing of the active list, if a new element arrives, it is linked to the waiting list and will participate the scheduling only after the current cycle of the above formula is done. Each time an element is moved to the waiting list, the weight count will be set to the original weight. See the example in Figure 379. This scheduling scheme is used to schedule the QOSs per port.

Smoothing the bursts

The WRR algorithm can cause long bursts of cells for the same FlowID/QOS if the weight is high. In order to make smaller bursts of data for the same FlowID/QOS, the QOSs are rotating in the active list several times before exhausting their weight. A Factor parameter is introduced per port, the weight of all QOSs is divided by the Factor and the result is the number of times they will be scheduled before moving to the end of the active list. For an example see Figure 380.

Negative Weight

When scheduling packets, the QOS of the specific FlowID will be scheduled until an EOP indication arrives, it is possible that the packet is bigger then the weight of the

QOS so that the count can become negative. Once EOF arrives if the weight count is still positive, the QOS can be scheduled for another packet. If the weight count is negative, the QOS is linked to the waiting list. The weight of the QOS is added to the weight count as well. Once the waiting list becomes active and the QOS is to be scheduled, if the weight count is still negative weight will be added again and the QOS will be sent to the waiting list again without scheduling.

This is similar to a negative credit, if a QOS consumes all the credits and beyond it has to accumulate enough credits in order to start scheduling again.

Best effort QOS

This is another flavor for the selection of the QOS per port. A dedicated QOS (QOS7) is called the best effort QOS. It does not participate in the strict priority scheme or in the weighted-round-robin scheme. This QOS is scheduled only if all QOSs that belong to the strict priority are empty and all the QOSs that belong to the WRR are empty as well. A simple way to describe it is that this QOS is scheduled only if all other QOSs are empty.

Strict Priority

Select a QOS can be done either by WRR algorithm or by assigning priorities to QOSs. Each port can assign priorities and WRR QOSs as shown in Figure 381.

Selecting QOS

As shown in Figure 382 there will be two QOS queues maintained per port. One is active list and other is waiting list. When a flow ID arrives from DBS, it will always be put in waiting list. Since there are eight QOSs per port each of this list is 8-bit wide and they are called QA_EMPTY(active list) and QW_EMPTY(waiting list). When corresponding bit is 1 then this QOS is empty and nothing to schedule. When QA_EMPTY is all 1s(no QOS to schedule) QW_EMPTY is copied to QA_EMPTY, and corresponding bits QW_EMPTY are made 1. Copying is done only for round robin priority QOSs. Strict priority QOSs will always remain in waiting list. At any point, when the port is to be serviced, based on the priority allocation for that port, first checked to find out if there are any strict priority QOSs are not empty in waiting list. If so then these QOSs are served first. Only when all strict priority QOSs are empty, then round robin QOSs are selected from active list. At any time if eop is sampled for the current QOS, this QOS will be served till eop irrespective of priority. Following algorithm explains how to select QOS for servicing once port is selected.

```

if (prev_qos_v == valid) then current_qos = prev_qos
else if (priority[0] & ~qw_empty[0]) current_qos = 0
else if (priority[1] & ~qw_empty[1]) current_qos = 1
else if (priority[2] & ~qw_empty[2]) current_qos = 2
else if (priority[3] & ~qw_empty[3]) current_qos = 3
else if (priority[4] & ~qw_empty[4]) current_qos = 4
else if (priority[5] & ~qw_empty[5]) current_qos = 5
else if (priority[6] & ~qw_empty[6]) current_qos = 6
else if (priority[7] & ~qw_empty[7]) current_qos = 7
else if (weight_mf != 0 & qa_empty[active_ptr] = 0)
current_qos = active_ptr

```

```

1  else if (&qa_empty[7:0]=0)
2      // repeat these steps for 0 to 7
3      case (prev_qos) // synopsys parallel_case
4      000:
5          case(1'b0)
6              qa_empty[1] : current_qos = 1
7              qa_empty[2] : current_qos = 2
8              qa_empty[3] : current_qos = 3
9              qa_empty[4] : current_qos = 4
10             qa_empty[5] : current_qos = 5
11             qa_empty[6] : current_qos = 6
12             qa_empty[7] : current_qos = 7
13             qa_empty[0] : current_qos = 0

```

To load pending QOSs from waiting to active. This function is performed any time active list is empty. Note that only non_strict priority QOSs are loaded to active list and strict priority QOSs are always serviced from waiting list. To load non priority QOSs to active list use

```

19  if (&qa_empty[7:0] = 1)
20      qa_empty[0] = priority[0] | (~priority[0] & qw_empty)
21      qw_empty[0] = ~priority[0] | (priority[0] & qw_empty[0]);

```

Levels of scheduling

Once a port is selected according to the port calendar, a QOS is selected between all the QOSs that belong to the port. Once a QOS is selected a FlowID is selected between the FlowIDs that belong to that QOS. See Figure 382.

Weight Quota

The weight of a QOS is calculated based on two values. The first value is the weight parameter, which is stored per QOS in the QOS descriptor memory and the second parameter is the weight_quota, which is a programmable value per all the QOSs. The weight of a QOS is the multiplication of the weight and the weight_quota.

Input Algorithm

1. Get FlowID and QOS from the shaper/DBS.
2. Read QOS memory:
 - $F_WP = QoS_mem[F_QOS][F_WP]$
 - $F_EMPTY = QoS_mem[F_QOS][F_EMPTY]$
 - $Q_NUM = QoS_mem[F_QOS][Q_NUM]$
 - $PORT_ID = QoS_mem[F_QOS][PORT_ID]$
3. Read Port memory:
 - $QA_EMPTY, QW_EMPTY = port_mem[Port_ID][QA_EMPTY, QW_EMPTY]$
 - $PRIORITY = port_mem[Port_ID][PRIORITY]$
4. Write to memories:
 - // Link FlowID to previous one
 - $Write\ flow_mem[F_WP][F_NEXT] \leq FlowID$
 - // Update tail of the link list of the FlowIDs
 - $Write\ QoS_mem[F_QOS][F_EMPTY] \leq 0$

```

1      • Write QoS_mem[F_QOS][F_WP] <= FlowID
2      • If (F_EMPTY == 1) Write QoS_mem[F_QOS][F_RP] <= FlowID
3      // Mark the QOS not empty in the port if needed
4      If (QA_EMPTY[QOS_NUM]==1)&(QW_EMPTY[QOS_NUM] == 1)
5          • Write port_mem[PORT_ID][QW_EMPTY[QOS_NUM]] <= 0
6      // Link port to calendar in case of Ingress chip configuration
7      If ((MODE==1) & (|QA_EMPTY) & (|QW_EMPTY)) begin
8          Write port_mem[INGRESS_PTR[15:8]] <= PORT_ID
9          INGRESS_PTR[15:8] <= PORT_ID
10         If(INGRESS_PTR[16]==0) INGRESS_PTR[7:0] <= PORT_ID
11         INGRES_PTR[16] = 0
12     end
13

```

Output Algorithm

The algorithm describes here doesn't show the conditions for selecting one of the options in step 9.

1. Get PortID from Port Calendar.
2. Read from port memory.
 - All parameters
3. If (| (QOS_EMPTY)) == 0) terminate output stage
Else continue to step 4
4. If active list is empty -> Move waiting list to Active list and clear waiting list.
5. Select QOS according to priority, WRR and previous QOS in middle of packet indication.
6. Read from QOS memory, address is {PortID, QOS}.
 - All parameters
7. Send to PFQ the F_RP value (This is the scheduled FlowID)
8. Read FlowID memory
 - All Parameters
9. Get port status from Status block.
10. Select one of the following options depending on parameters from all memories:
 - o delink_flowid
 - o move_flowid_to_end_of_active
 - o don't_move_flowid_and_update_cell_cnt
 - o do_nothing_to_flowid
 - o delink_QoS
 - o move_QoS_to_end_of_active
 - o don't_move_QoS_and_update_weight_m
 - o move_QoS_to_waiting
 - o do_nothing_to_QoS
 - o For Ingress chip only:
 - move_port_to_end_of_link_list
 - delink_port
11. Update internal and external memories according to above decision.

Algorithm to calculate q_weight_m during output phase

```

47  if (qos_empty)
48      q_weight_m <= 25'h0
49  else if (q_weight_m ≤ 0 & eop)           // move to waiting list
50      q_weight_m <= q_weight_m - 1 + (q_weight * weight_quota)
51  else if (q_weight_m ≠ 25'h1000000)      // not the maximum negative value
52      q_weight_m <= q_weight_m - 1

```

Memory links for Input Phase

The input phase starts with a FlowID from the DBS. With this parameter all memories are accessed in sequence to get all relevant parameters required for the input phase. See Figure 383.

Memory links for Output Phase

The output phase starts with a port assigned by the port calendar inside the DBS block. The port parameters are read from the internal memory, and a QOS is selected for that port based on the none-empty QOSs and the strict priority. The address to the QOS memory is {PortID, QOS_SELECTED}. The read pointer in the QOS parameters is the FlowID to send to the PFQ. Once reading the parameters of that FlowID it can be sent to the PFQ together with the EOP indication. See Figures 384 and 385.

CPU Port

CPU port is different then the rest of the 64 ports supported. The CPU port does not have any QOS and it's rate is fixed to OC3. FlowIDs are linked to the CPU ports using the external NEXT memory as the next field. FlowIDs are scheduled based on a simple round robin between them.

Blocks Description

The scheduler is divided to several different functional blocks as follow:

Input Phase Block

The input phase contains a small FIFO to accumulate the link requests from the shaper. These link requests are served once an input slot is available. During this stage a FlowID is linked in a link list and if needed a QOS is linked to the waiting lists. Also the port associated with that FlowID is marked as non-empty.

Output Phase Block

The output phase block receives a port from the port calendar. Based on the status of that port, a decision is made if to go ahead with scheduling (if the port is not full) or to dedicate the bandwidth to the software (if the port is full). If scheduling is ok, a QOS is selected according to a weighted-round-robin algorithm and based on that QOS a FlowID is selected. This FlowID is then transferred to the PFQ.

If the cell count of a FlowID becomes zero, means a full packet was sent, the FlowID would be de-linked from the link list if there is no other packets linked to this FlowID. It is possible that because of that the QOS will be de-linked as well.

Software Interface

Software needs to access the scheduler database for several purposes:

1. Initialization after reset: port calendar and QOS databases as well as internal registers.
2. When setting-up a connection, to write the FlowID parameters.
3. For testing purposes of the internal memories mainly.

All software accesses to the scheduler are done through specific registers during the dedicated time slots.

Memory Description

The following is a description of the memories needed to implement the scheduling algorithm.

External Memory – Next FlowID

Address: Flow ID, Length: 1M (20 bits wide)

See Figure 386.

Internal memory – QOS Parameters

Address: QOS, Length: 1K (17 bits wide)

See Figure 387.

Internal memory – QOS Descriptors

Address: QOS, Length: 1K (66 bits wide)

See Figure 388.

Internal memory – Port Parameters

Address: Port ID, Length: 64 (21 bits wide)

See Figure 389.

Internal memory – Port descriptor

Address: Port ID, Length: 64 (47 bits)

See Figure 390.

Internal memory – Next Port

Address: Port ID, Length: 64 (8 bits)

See Figure 391.

Software Interface

The scheduler can be programmed via the software interface by a CPU. All the scheduler's control registers and memories can also be written and read for test purposes through the software interface. The following sections list all the control registers and all the scheduler software interface commands.

Registers

See Figure 392.

Commands

Read FID NEXT Memory

See Figure 393.

- 1 **Write FID NEXT Memory**
- 2 See Figure 394.
- 3 **Read QOS Parameters**
- 4 See Figure 395.
- 5 **Write QOS Parameters**
- 6 See Figure 396.
- 7 **Read QOS Descriptor**
- 8 See Figure 397.
- 9 **Write QOS Descriptor**
- 10 See Figure 398.
- 11 **Read PORT Parameters**
- 12 See Figure 399.
- 13 **Write PORT Parameters**
- 14 See Figure 400.
- 15 **Read PORT Descriptor**
- 16 See Figure 401.
- 17 PQ[0] = PREV_QOS_VALID
- 18 PQ[3:1] = PREV_QOS[2:0]
- 19 PQ[6:4] = ACTIVE_PTR[2:0]
- 20
- 21 **Write PORT Descriptor**
- 22 See Figure 402.
- 23 **Read PORT Next**
- 24 See Figure 403.

1 **Write PORT Next**

2 See Figure 404.

3 **Init QOS**

4 See Figure 405.

5 This command write the parameters to the QOS Parameter Memory and also
6 initiates a write to the QOS Descriptor Memory with the following values:

7 F_EMPTY = 1

8 Rest of the fields should be reset to zero

9 **Init port**

10 See Figure 406.

11 This command write the parameters to the Port Parameter Memory and also
12 initiates a write to the Port Descriptor Memory with the following values:

13 QW_EMPTY[7:0] and QA_EMPTY[7:0] = 8'b1111_1111

14 Rest of the fields are reset to zero.

15 **Reserved Op-codesReserved Op-codesReserved Op-codes**

16 The op-codes 1111 are reserved.

17 **Timing Diagrams**

18 **Global Timing for Input/Output Stages**

19 The input phase and the output phase should be combined during the 8 cycles
20 budget per cell. Since we do so, there are data dependencies between input phase and
21 output phase. Figure 407 shows the memory accesses needed for each memory for every
22 8 cycles of input/output phase (Figure 408 shows the addresses, all data is 2 cycles later
23 for the external memory):

24 **Input/Output stage Pipes**

25 **Input Stage**

26 The input stage actually starts one cycle before cycle1 by sampling a valid
27 signal from the DBS that holds for the entire input stage cycles. If there is no valid signal
28 asserted, the input pipe is frozen by not changing the first pipe stage (this will also save
29 power). See Figure 408.

30 **Output Stage**

31 When working as an egress chip, the output stage starts with reading the port
32 memory parameters using a port ID sent by the port calendar inside the DBS. When in an

ingress chip mode the location to read from is determined by the INGRESS_PTR, which is a register controlled by the scheduler.

The output stage diagram doesn't show the data dependencies to the input phase. Data dependencies will be discussed in the next section. See Figure 409.

Data Dependencies

Data dependencies between the output stage and the input stage

In the scheduler, the data dependencies might happen when the input stage and the output stage access to the same memory address. If one stage reads a memory location before the other finishes updating the same memory location, then the read result might be stale. Or if both stages write the same memory address, the later write might overwrite the previous write result.

There is no data dependency issue for the FlowID memory since the input and output stages are mutually exclusive on a specific FlowID, i.e. the input and output stages never work on the same FlowID. If the input stage is linking a FlowID into the linked list, the output stage is not aware of the existence of the FlowID yet so there is not output phase for this FlowID. If the output stage is scheduling a FlowID, the input stage can not link the same FlowID by the principle of operation. Even if by looking at the FID Next memory operations, there is no overlapping between the input stage and the output stage memory accesses.

For the QOS memory accesses, the input stage only updates the field of F_WP at cycle 7. But the output stage may also update the F_WP field at cycle 6 of the next 8 cycle slot. To resolve the data dependency issue, when the QOS address is the same for the input and the output stages, the input stage updates not only the memory field of F_WP at cycle 7 but also the register field of F_WP which stores the output stage read result of cycle 3.

By the same token in the port memory case, when the input and the output stages access to the same port memory address, the input stage updates the memory field of QW_EMPTY at cycle 7 and at the same time it has to update the same QW_EMPTY field of a register which stores the output stage read result at cycle 2 to avoid the data dependency.

The output stage data dependencies between consecutive 8 cycle time slots

The port memory accesses of the output stage take more than 8 cycles, which implies that the port memory accesses in one time slot overlaps the port memory accesses in the following time slot. More specifically, before the previous 8 cycle time slot updates the port memory at cycle 6, the current 8 cycle slot has to read the port memory at cycle 5.

To avoid the data dependencies when the consecutive 8 cycle time slots access to the same port (there is only one port in this case), the read result has to be taken from a register which stores the write date of the previous slot memory write instead of the memory content.

1 Traffic Shaper/Meter

2 This block can operate in 2 modes:

- 3 • Traffic Shaper mode: in this mode the Shaper shapes specific rates to FlowIDs
4 according to Single or Dual Leaky Bucket algorithm. The shaper selects a Peak rate
5 or a Sustained rate per FlowID depending on an accumulated credit. Once a FlowID
6 has been shaped, it is marked as such in a common.
- 7 • Traffic Meter mode: in this mode the Meter measures the rate of the incoming
8 FlowID traffic and mark the FlowID as a violator if the rate is above a specific
9 threshold which is programmed for a FlowID. The user periodically read the
10 indications in the FlowID memory.
11 See Figure 410.

13 4.1 Features

- 14 1. This block operates at 200MHz clock
- 15 2. A FlowID can be issued from the Database at max once every 8 cycles of
16 200MHz clock.
- 17 3. A FlowID is driven to the Database at max once every 8 cycles of 200MHz clock.
- 18 4. Shapes input traffic per FlowID using Single or Dual-Leaky-Bucket algorithm.
- 19 5. Shapes traffic for rates from ~48.8kbps up to 9.920Gbps.
- 20 6. Meter traffic from 12.8kbps up to 9.920Gbps in increments of 64kbps.
- 21 7. Up to 1024 different programmable rates can be shaped simultaneously. Each
22 shaped FlowID will point to one of the entries in this table for a specific shaping
23 rate.
- 24 8. Up to 1024 different programmable thresholds can be metered simultaneously.
25 Each metered FlowID will point to one of the entries in this table for a specific
26 threshold rate.
- 27 9. Each FlowID is programmable to have Dual Leaky Bucket, Single Leaky Bucket
28 or no shaping at all.
- 29 10. Can shape or meter up to 1M FlowIDs simultaneously.
- 30 11. Supports microprocessor programming of the memories for FID, rates, etc.
- 31 12. 28 bits of the external memory in FID1 are used for dual leaky bucket algorithm
32 only. It is possible not to use these bits if using single leaky bucket only. Note that
33 for a Meter function the extra memory for the dual-leaky-bucket is not needed.
- 34 See Figure 411 and 412. Note that the memory bits, which are not used for the single
35 leaky bucket should be tied to zero.
- 36 See Figure 413. This block doesn't need all the memories for every single mode of
37 operation. Refer to the memory description chapter to see the usage of the memories in
38 each mode.

40 Functional Description

41 The traffic shaper receives an input phase from the Database block along with the
42 RateID and then starts shaping the FlowID on a per cell basis. Each time a cell is shaped
43 the Shaper generates a message to the Database. The Database then responds with an
44 empty indication. All packets link list and output to the PFQ are done in the Database
45 block.

Leaky-Bucket Algorithm

The basic leaky-bucket algorithm is an accumulation of credits for the given FlowID, once the amount of credit is enough for a cell, this cell is considered shaped and is transferred. The equation to calculate is:

$$R_{out} = R_{in} * M/N$$

When R_{in} is the input rate (based on the selected class), M input credit, N output credit. The following Figure shows an example of $M/N = 2/5$:

See Figure 414

Dual-Leaky-Bucket Algorithm

A dual-leaky-bucket contains 2 virtual buckets per FlowID (virtual because credit is accumulated only in one bucket). One bucket shapes the traffic for a sustained data rate and the other shapes the traffic for a peak data rate.

The shaped traffic starts by using the M_s/N_s values of the sustained data rate. If for some reason the accumulated credits are above a specific threshold amount, then the shaped traffic uses the M_p/N_p values of the peak data rate until the amount of credits returns to zero.

The credit is calculated using the current time measurement. When shaping a cell, the difference between the last shaping time and the current shaping is time calculated to create ΔT this is the elapsed time from the last shaped cell. $\Delta Credit$ is calculated as follow:

$$\Delta Credit = \Delta T * M_s - N_s$$

$$New_Credit = Old_Credit + \Delta Credit$$

If ΔT is long enough then a positive credit is accumulated. If ΔT is short then a negative credit is accumulated. Once **New_Credit** is above the threshold then the shaping is done with the peak rate while the credit is always calculated with the sustain data rate.

Implementing Timing Wheel for leaky-bucket algorithm

The timing wheel contains 256K time slots. The wheel is rotating in a rate of one time slot every cell cycle. Each time slot can have few FlowIDs link to it. There is one exit point in the timing wheel, which is time zero. A FlowID is linked to a future time slot based on its shaping parameters. Once that time slot reached time zero, all the FlowIDs that are linked to that specific time slot are dropped to an output FIFO. When a FlowID enters the output FIFO it means that one cell for that specific FlowID is shaped. Once every cell cycle one FlowID is taken from the output FIFO and is being serviced, means, a new time slot is calculated for the FlowID based on the rateID and some other parameters. The FlowID then is linked to that future time slot. In parallel the shaper sends this FlowID to the Database block, in this block a cell count and packets link list are maintained per FlowID.

It is possible that several FlowIDs are linked to the same time slot. When this time slot arrives to time zero, several FlowIDs are scheduled to the output FIFO at the same cell cycle. In this case the wheel services only one FlowID when the rest are waiting in an output FIFO to be serviced in the next cell cycle. See Figure 415.

The timing wheel is actually fixed and only the zero reference point is moving. That means the wheel can be implemented as a memory and the zero point is a moving address, which is incremented once every cell cycle.

A single leaky bucket algorithm is implemented using the timing wheel. As described above, the Equation to implement the Leaky Bucket is:

$$R_{out} = R_{in} * M/N$$

When assigning K to N/M, we will get

$$R_{out} = R_{in} / K$$

The K parameter is the time in the future to place the FlowID. For example if K=3 means the R_{out} is 1/3 of R_{in} then the FlowID is linked to the time slot which it's distance from the current time is 3. Since the timing wheel is rotating once every cell cycle (once every 8 cycles) the next time the FlowID will be seen is after 3 cell cycles when the time slot will reach zero:

See Figure 416.

Fractional value of K

It is possible that K will have a fractional value and an integer value for example to get 40% of the bandwidth K should have the value of 2.5. It is impossible to have a fractional time slot, therefore, to place a FlowID in a future time slot only the integer part of K is used. The fractional part of K is getting accumulated and is added to the integer part only if it is greater or equal to 1. For example, if K=1.28 (78.125% of the line rate) the shaping will appear as in Figure 417.

Output FIFO in the timing wheel

The output FIFO is a link list of the FlowIDs that are shaped for one cell and are waiting to be re-linked to the wheel to some future timing slot. The output FIFO is used because it is possible that more than one FlowID will be linked to a specific slot. In that case, when the slot arrives to time 0, several FlowIDs are shaped at the same time. Since the Shaper can handle only one FlowID every cell cycle, the rest of the FlowIDs have to wait.

Each time a new list of FlowIDs are linked to the output FIFO, they are marked with the time arrival to the output FIFO. When calculating the future time, the result is decreased by the amount of time the FlowID was in the output FIFO. For example, 2 FlowIDs are linked in the output FIFO, $K_{FID1}=3$ and $K_{FID2}=5$. The scheduled time slots will appear as in Figure 418.

Dual-Leaky-Bucket using timing wheel

To create a dual-leaky-bucket in timing wheel there is a need for the following parameters:

- K_s for sustain data rate.
- K_p for the peak data rate.
- M_s, M_p Sustain and Peak input rates.
- **Threshold** the threshold value of the accumulated credits to switch to peak rate.
- **Credit** the total amount of positive credits (always a positive number).
- **Residue** the accumulation of the fractional part of K_s or K_p .
- **Time_A** the time of arrival to the output FIFO.

- 1 • **Time_L** time of last scheduled cell to a time slot.
- 2 • **Back Ground Process** is used to mark timer wrapping around since measuring ΔT .
- 3 There are 3 processes for a FlowID: input phase, output phase and a background process:
- 4 • Input phase is when the FlowID is first linked to the timing wheel, same calculation
- 5 as the output phase is applied except the residue in this case is 0.
- 6 • Output phase is when the FlowID is services in the output FIFO.
- 7 • Background process is a round robin service between all FlowIDs (1M of them) to
- 8 mark the wrapping around of the current time. This is done for a FlowID that
- 9 becomes empty and after a while becomes non empty, the credits should be calculated
- 10 properly for the FlowID.
- 11 When changing the sustain rate to a peak rate then the peak rate is used until the
- 12 amount of credits is down to zero. The calculation of $\Delta Credit$ is being done using the
- 13 sustain data rate always. Figure 419 describes the use of sustain and peak rates.

15 Calculation of $\Delta Credit$

16 The calculation of the $\Delta Credit$ depends of the rates and the time elapsed from the
 17 last output phase. There is a small difference between the calculations in the output or the
 18 input phases.

19 **K** is a pure number. **K** multiple by one time tick gives the amount of time ticks
 20 needed to shape one cell. The equation is as follow:

$$21 \quad \Delta Credit = (\Delta T - 1 \text{ TimeTick} * K) / K = (\Delta T - K) / K$$

22 ΔT is the time elapsed from the last output phase. **K** is decremented from ΔT since
 23 during that time a cell was shaped. The division by **K** allows measuring the credits in
 24 terms of cells. Note that the **K** in the division is the pure number (which is N/M) when
 25 the **K** inside the parenthesis has units of time.

- 26 • When the output phase uses the sustain or peak rates the equation to be used is:

$$27 \quad \Delta Credit = (\Delta T - K_s) / K_s = \Delta T / K_s - 1$$

- 28 • When the input phase calculates the change in the credit there is no cell shaped
 29 therefore:

$$30 \quad \Delta Credit = \Delta T / K_s$$

31 Equations (Algorithm) for Input phase

- 32 1. $NewCredit = 0$
- 33 2. $TimeSlot = INT(K_s) // \text{integer part of } K_s$
- 34 3. $NewResidue = FRAC(K_s) // \text{fractional part of } K_s$

35 Equations (Algorithm) for Output phase

- 36 1. $\Delta T = CurrentTime - Time_L$
- 37 2. Calculate $\Delta Credit$
- 38 3. $TempCredit = Credit + \Delta Credit$
- 39 4. $If(TempCredit > Threshold)$

```

1      NewCredit = Threshold
2      Select K = Kp
3      Else
4      NewCredit = TempCredit
5      Select K = Ks
6      5. NewTime = (INT(K + Residue) - TimeA
7      6. TimeSlot = if(NewTime < 0) ? 0 else NewTime
8      7. NewResidue = FRAC(K + Residue)
9      8. NewMarkBits = 11

```

10 Equations (Algorithm) for Background phase

```

11      1. if(MarkBits == 11) NewMarkBits = 10
12      2. else if(MarkBits == 10) NewMarkBits = 0
13

```

14 **R_{in} and K factor**

15 During input phase or output phase a future time slot is calculated according to K
16 and time variables. A slow rate, means a high K value, will cause a linking of the FlowID
17 to a far timing slot. Whereas a high rate, means a low K value, will cause a linking of a
18 FlowID to a closer timing slot. This chapter is to describe the exact R_{in} for the calculation
19 of:

$$20 \quad R_{out} = R_{in} / K$$

21 The chip runs at 200MHz and the data path is 64 bits, therefore,

$$22 \quad R_{in} = 200 \text{ MHz} * 64 \text{ bits} = 12.8 \text{ Gbps}$$

23 There are 256K timing slots in the timing wheel, therefore, the slowest rate that it
24 can process is:

$$25 \quad R_{out}(\text{min}) = 12.8\text{G}/256\text{K} = 48.828125 \text{ kbps}$$

26 In order to have a shaped rate of 64kbps, K should be 200,000:

$$27 \quad 64\text{kbps} = 12.8\text{G}/200000$$

28 That means that if a FlowID is placed in the 200000th time slot in the future each
29 time, its output rate is exactly 64kbps. Figure 420 shows some of the standard rates with
30 their corresponding K. Note that the ratio between OC192 (9.92G) to 12.8G is 1.290

32 **Violator indication**

33 When using a dual-leaky-bucket algorithm the violator indication should be set
34 when using the peak rate. When using a single-leaky-bucket algorithm, the violator
35 indication should not be set.

37 **Metering**

38 Metering is the other mode of operation of this block. The user can assign one of
39 1K different rates pre defined by the user to a specific FlowID and the block will monitor
40 the traffic to that FlowID. If the actual traffic rate of the FlowID, the measured rate, is
41 greater then the assigned rate then the FlowID is marked by the hardware.
42 The user should periodically read the marking bit per active FlowID. Once the user reads
43 the FlowID information, the bit is cleared and remains clear until the input traffic is
44 violating again.

External memory for Metering mode

The meter block used the same memory bits of the FlowID but for different meaning. Each FlowID needs to have a counter to count number of incoming cells. Also each FlowID needs to have a bit to mark the FlowID as a violator of the assigned rate or not.

Background process of Metering

The background process of the metering is used to check for the accumulated amount of cells per FlowID every specific amount of time. Once every 8 cycles the background process services one FlowID. That means, a FlowID is being checked by the background process once every 40mSec (at 200MHz clock).

To provide more accuracy to lower rates, a counter is added to each FlowID to control the amount of times the FlowID is accessed by the background process before the actual measurement and marking is done. When the background process accesses a FlowID and this counter is greater than zero, then the background process only decrement that counter, when the counter is zero the background process compares the accumulated amount of cells to the threshold, mark the FlowID if the threshold is smaller than the accumulated amount and then the cell counter is reset, the background counter is set to the value assigned by the user.

Once the metering algorithm marked a FlowID as a violator, this FlowID is to remain marked until the user reads the FlowID memory. When the user reads the FlowID memory, the mark bit is reset by the hardware.

Output FIFO for Metering

During the input phase of a FlowID in the metering mode, the FlowID can be marked as a violator. If it is to be marked as a violator, then the FlowID is linked to a link list of violators. The software accesses this list when it wants to get the list of the violators, this way the software doesn't have to read all FlowIDs to detect the violators. Once the FlowID is linked to that list, only when the software reads it, it will be removed from the list. The software first reads the value of the MARK_RP register, which is the pointer to the head of the violators link list and then reads the FlowID memory. The Meter generates a write to the FlowID memory to clear the mark bit,

Block Description

Shaper's Input Phase

The Shaper input phase receives a rateID, FlowID and valid signal from the DataBase block. This is a request from the database block to shape the FlowID. Once the FlowID is inside the shaper, no more input phase is initiated by the DBS block. If the specific time slot was empty then it is marked as non-empty, if the specific time slot was not empty then the FlowID is linked to the last FlowID in the link list.

Shaper's Slot Background process

The Traffic Shaper Wheel is constantly rotating one time slot every 8 clocks of the 200MHz clock. As each time slot becomes current, the background process is activated to get all the FlowIDs from the current time slot and link them to the output

FIFO link list. While linking the FlowIDs from the time slot to the output FIFO, the time slot is marked as empty and become the last time position in the wheel.

Shaper's FlowID Background process

This background process is to mark the time in each FlowID such that in the next output phase a ΔT is calculated correctly even if the free running time counter was wrapped around. The process resets the bits while the output phase sets the bits. Each FlowID is serviced once every 1M cell cycles (8M clock cycles). In order to mark a FlowID to use only the threshold value as a credit without using ΔT , this process has to be repeated 2 times. Therefore, a FlowID is marked as "too old to calculate credit" at least after 1M+1 cell cycles and no more then 2M cell cycles (16M clock cycles). In this case the credit value becomes the value of the threshold.

Shaper's Output Phase

The output phase operates once every cell cycle on the FlowID, which is the head of the output FIFO's link list. Operating on a specific FlowID means shaping one cell of the FlowID. The shaper drives the FlowID, which is the head of the output FIFO's link list to the database block. The database block manages the link list of the packets pr FlowID. The act of driving a FlowID to the database block means that one cell is shaped for that FlowID.

The database can return a message to the shaper that the current FlowID is empty, in this case the shaper will not link the FlowID to a future time slot. Note that the Shaper always shapes a cell at a time for packets traffic and for cells traffic. It is the Database block that keeps track of shaped packets/cells and generate de-queue commands to the PFQ accordingly. During metering mode, the output phase is disabled.

Shaper's Time measurement

For the dual-leaky-bucket algorithm there is a need to measure time for calculating the change in the credit. A 21 bits register is used to count the time, this counter is incremented once every 8 cycles. For all calculations the logic used only 18 bits of the time counter except for the FlowID background process.

Since there are 1M FlowIDs supported there is a need to measure time up to 1M units and more, that's is why the time counter has 21 bits.

Traffic Shaper Algorithms

The below algorithms are implementation specific and not only the mathematical equations and conditions. To understand the algorithms, one should look at the memory description first. Memory description.

Shaper Input Phase Algorithm (for single leaky bucket)

1. Receive from database block: *VALID, FID, RATEID*
2. *RD shp_rate_id_mem[RATEID]{Ks, Kp, Threshold, M}*
3. *RD FID2_mem[FID]{start, bg_stt, time_msb}*
RD FID1_mem[FID]{last_time, residue, credit, p_s}
4. // Calculate future time slot and residue based on previous values
if((p_s==1) & bg_stt[1:0]!=2'b00)

```

1      a) slot_id[17:0] <= cur_slot_addr[17:0] + INT(Kp)
2      b) n_residue[17:0] <= FRAC(Kp)
3      c) if(VALID & INT(Kp)==0) link_to_out_fifo <= 1 else link_to_out_fifo <=
4          0
5      else
6      a) slot_id[17:0] <= cur_slot_addr[17:0] + INT(Ks)
7      b) n_residue[17:0] <= FRAC(Ks)
8      c) if(VALID & INT(Ks)==0) link_to_out_fifo <= 1 else link_to_out_fifo <=
9          0
10     5. RD slot_mem[slot_id]{rp_slot, wp_slot, e_slot}
11     6. // Calculate new credit, see block diagram below
12         n_credit = f(current_time, last_time, Ks, M, credit, threshold, bg_stt)
13     7. // Calculate rate selector for next output phase
14         if(n_credit>=threshold) n_p_s <= 1
15         else if(n_credit==0) n_p_s <= 0
16         n_p_s <= p_s
17     8. // Create link FID to future time slot
18     9. if(e_slot==1)
19         a) {n_rp_slot, n_wp_slot, n_e_slot} <= FID, FID, 0
20         b) n_start <= 1
21     10. else
22         a) {n_rp_slot, n_wp_slot, n_e_slot} <= rp_slot, FID, 0
23         b) n_start <= 0
24     11. if(link_to_out_fifo)
25         Link FID to output FIFO
26     12. else if(VALID==1)
27         WR slot_mem[slot_id]{rp_slot, wp_slot, e_slot} <= n_rp_slot, n_wp_slot,
28             n_e_slot
29     13. // Write back for FID fields
30     if(VALID==1)
31         FID1_word_en[0], Wr FID1_mem[FID]{residue} <= n_residue
32         FID1_word_en[3], Wr FID1_mem[FID]{credit, p_s} <= n_credit, n_p_s
33         FID2_word_en[0], Wr FID2_mem[FID]{start} <= n_start
34     if(VALID==1 & e_slot==0) // to avoid a write to an unknown FID's next field
35         FID2_word_en[1], WR FID2_mem[wp_slot]{next_FID, bg_stt,
36             time_msb} <= FID, 11, cur_time_reg[20]
37 Slot background Phase Algorithm
38     1. current_slot <= current_slot + 1; //wraps around 256k
39     2. RD shp_slot_mem[current_slot]{rp_slot, wp_slot, e_slot}
40     3. WR shp_slot_mem[current_slot]{e_slot} <= 1'b1
41     4. If(e_slot==1)
42         a. do nothing //nothing to add to the output registers
43     5. else
44         a. wp_out_reg <= wp_slot
45         b. wr_en_shp_FID_next_mem <= (e_out_reg==0)
46         c. if(e_out_reg==0) //output fifo is not empty

```

```

1          i. rp_out_reg <= rp_out_reg
2          ii. e_out_reg <= e_out_reg
3      d. else
4          i. rp_out_reg <= rp_slot
5          ii. e_out_reg <= 0
6      6. If (wr_en_shp_FID_next_mem==1)
7          FID2_word_en[0], WR shp_FID_next_mem[wp_out_reg]{next_FID} <= rp_slot
8      7. If (e_slot==0)
9          FID1_word_en[1], WR FID1_mem[rp_slot]{arrival_time} <= current_time

```

11 FID Background Phase Algorithm (for dual-leaky-bucket only)

```

12      1. current_FID <= current_FID + 1;
13      2. RD FID2_mem[current_FID]{empty, bg_stat, time_msb, start}
14      3. case (bg_stat) //synopsys full case parallel case
15          a. 2'b11: n_bg_stat <= 2'b10
16          b. 2'b10: n_bg_stat <= 2'b00
17          c. default: 2'b00
18      endcase
19      4. FID2_word_en[1], WR FID2_mem[current_FID]{empty, bg_stat, time_msb, start} <=
20          n_empty, n_bg_stat, time_msb, start

```

21 Shaper Output Phase Algorithm

```

22      1. RD FID1_mem[rp_out_reg]{residue, arrival_time, last_time, credit, p_s}
23      2. RD FID2_mem[rp_out_reg]{next_FID, bg_stt, time_msb, start}
24      3. Send to Database block: rp_out (as O_FID), ~e_out_reg (as O_VALID)
25      4. Receive RateID from Database
26      5. RD rateid_mem[RateID]{Ks, Kp, threshold}
27      6. // Calculate future time slot and residue based on previous values
28          time_spent_in_fifo[17:0] <= current_time[17:0] - arrival_time[17:0]
29          if (p_s==1)
30              a) slot_id[18:0] <= cur_slot_addr[17:0] + INT(Kp) -
31                  time_spent_in_fifo[17:0]
32              b) n_residue[17:0] <= FRAC(Kp)
33              c) if (VALID & INT(Kp)==0) link_to_out_fifo <= 1 else link_to_out_fifo <=
34                  0
35          else
36              d) slot_id[18:0] <= cur_slot_addr[17:0] + INT(Ks) -
37                  time_spent_in_fifo[17:0]
38              e) n_residue[17:0] <= FRAC(Ks)
39              a. if (VALID & INT(Ks)==0) link_to_out_fifo <= 1 else link_to_out_fifo <=
40                  0
41      7. // Calculate new credit
42          a. n_credit = f(current_time, last_time, residue, Ks, Kp, M, credit, threshold)
43      8. // Calculate rate selector for next output phase
44          if (n_credit >= threshold) n_p_s <= 1
45          else if (n_credit == 0) n_p_s <= 0
46          n_p_s <= p_s

```

```

1  9. // prepare values for FID link to a slot
2      RD slot_mem[slot_id][rp_slot, wp_slot, e_slot]
3      if (e_slot==1)
4          a. {n_rp_slot, n_wp_slot, n_e_slot} <= FID, FID, 0
5          b. n_start <= 1
6          c. next_wr_en <= 0
7      else
8          a. {n_rp_slot, n_wp_slot, n_e_slot} <= rp_slot, FID, 0
9          b. n_start <= 0
10         c. next_wr_en <= 1
11 10. Receive EMPTY indication from Database
12 11. // Link FID to future time slot
13     if (O_VALID==1 & ~EMPTY & link_to_output_fifo)
14         link FID to output FIFO
15     else if (O_VALID==1 & ~EMPTY )
16         WR slot_mem[slot_id][rp_slot, wp_slot, e_slot] <= n_rp_slot, n_wp_slot,
17         n_e_slot
18 12. // Write back all FID variables
19     if (O_VALID & ~EMPTY)
20         WR FID1_mem[rp_out_reg][residue] <= n_residue
21         WR FID1_mem[rp_out_reg][last_time] <= current_time_reg
22         WR FID1_mem[rp_out_reg][credit, p_s] <= n_credit, n_p_s
23         WR FID2_mem[rp_out_reg][bg_stt, time_msb, start] <= 11, cur_time_reg[20],
24         n_start
25         If(next_wr_en==1) WR FID2_mem[wp_slot][FID_next] <= rp_out_reg
26 13. if(rp_out_reg==wp_out_reg & don't place FID in output fifo again) e_out_reg
27     <= 1
28     else e_out_reg <= 0
29 14. rp_out_reg <= next_FID
30

```

Translate {Exp, Mantissa} presentation of K to Integer and Fraction

The K factors are represented in the database with mantissa and exponent.

- For Ks and Kp the presentation is {1.mantissa} x 2^{exp}. exp means the number of places to shift the decimal point to the right.
- For inverse of Ks the presentation is {0.mantissa} x 2^{-exp}. exp means the number of places to shift the decimal point to the left.
- Exp can be common for the Ks value and the InvKS value in the internal memory.

The following algorithm is to translate Ks or Kp value to integer and fraction values from a presentation of {exp, mantissa}. This code is used in the input phase and in the output phase. rate_int = INT(K) and rate_frac = FRAC(K):

```

1. case(exp) // synopsys full case parallel case
0: rate_int<={17'b0,1} rate_frac<={man(17:0)}
1: rate_int<={16'b0,1,man(17)} rate_frac<={man(16:0),1'b0}
2: rate_int<={15'b0,1,man(17:16)} rate_frac<={man(15:0),2'b0}
3: rate_int<={14'b0,1,man(17:15)} rate_frac<={man(14:0),3'b0}
4: rate_int<={13'b0,1,man(17:14)} rate_frac<={man(13:0),4'b0}

```

```

1      5: rate_int<={12'b0,1,man(17:13)} rate_frac<={man(12:0),5'b0}
2      6: rate_int<={11'b0,1,man(17:12)} rate_frac<={man(11:0),6'b0}
3      7: rate_int<={10'b0,1,man(17:11)} rate_frac<={man(10:0),7'b0}
4      8: rate_int<={9'b0,1,man(17:10)} rate_frac<={man(9:0),8'b0}
5      9: rate_int<={8'b0,1,man(17:9)} rate_frac<={man(8:0),9'b0}
6      10: rate_int<={7'b0,1,man(17:8)} rate_frac<={man(7:0),10'b0}
7      11: rate_int<={6'b0,1,man(17:7)} rate_frac<={man(6:0),11'b0}
8      12: rate_int<={5'b0,1,man(17:6)} rate_frac<={man(5:0),12'b0}
9      13: rate_int<={4'b0,1,man(17:5)} rate_frac<={man(4:0),13'b0}
10     14: rate_int<={3'b0,1,man(17:4)} rate_frac<={man(3:0),14'b0}
11     15: rate_int<={2'b0,1,man(17:3)} rate_frac<={man(2:0),15'b0}
12     16: rate_int<={1'b0,1,man(17:2)} rate_frac<={man(1:0),16'b0}
13     default: rate_int<={1,man(17:1)} rate_frac<={man(0),17'b0}
14     // make case 17 and above default to keep, for exp>=17 use 17 only
15     endcase

```

2. $INT(K) \leq rate_int - 1$ $FRAC(K) = rate_frac$

Calculation of K based on desired Rout

Since the shaper performs the following equation:

$$R_{out} = R_{in} / K$$

K value can easily be calculated as:

$$K = R_{in}/R_{out}$$

For rates units of [bit/Sec] K is calculated like this:

$$K = 12.8 * 10^9 / R_{out}$$

Credit calculation diagram

For input phase the credit calculation depends on the following parameters:

$n_credit = f(current_time, last_time, Ks, M, credit, threshold)$.

For output phase the credit calculation depends on the following parameters:

$n_credit = f(current_time, last_time, residue, Ks, Kp, M, credit, threshold)$.

Both calculations are almost the same, the implementation should be the same for both and only some inputs may change if using it for input phase or output phase.

Output phase calculation is shown in Figure 421.

Meter Algorithms

Input Algorithm

1. Receive FID, RATEID, VALID, NCELL from database block
2. **RD** rateid_mem[RATEID] {threshold}
3. **RD** FID1_mem[FID] {cell_cnt, mark, bg_cnt, bg_cnt_value}
4. $n_cell_cnt_tmp[27:0] = cell_cnt[26:0] + NCELL[15:0]$
5. if($n_cell_cnt_tmp[27] == 1$) $n_cell_cnt \leq 27'h7ffffff$
else $n_cell_cnt \leq n_cell_cnt_tmp[26:0]$
6. $n_mark \leq mark \mid (n_cell_cnt_tmp[27:0] > \{1'b0, threshold[26:0]\})$
7. $n_bg_cnt \leq bg_cnt$
8. if($mark=0$ & $n_mark==1$) // this FID just violated the metering now – link it to viol_reg
if ($MARK_EMPTY==1$)

```

1      {n_MARK_RP, n_MARK_WP, n_MARK_EMPTY} <= FID, FID, 0
2      wr_next_FID <= 0
3      else
4          {n_MARK_RP, n_MARK_WP, n_MARK_EMPTY} <= MARK_RP, FID, 0
5          wr_next_FID <= 1
6      else // either the FID did not violate metering or it is already marked
7          {n_MARK_RP, n_MARK_WP, n_MARK_EMPTY} <= MARK_RP,
8          MARK_WP, MARK_EMPTY
9          wr_next_FID <= 0
10     9. if(VALID==1)
11         WR FID1_mem[FID] <= {n_cell_cnt, n_mark, n_bg_cnt, bg_cnt_value}
12     10. if(VALID==1 & wr_next_FID==1)
13         WR FID2_mem[MARK_WP][next_FID] <= FID
14

```

Background Algorithm

```

15
16     1. Read FID_mem[FID] {cell_cnt, mark, bg_cnt, bg_cnt_value}
17     2. if(mark==1) //freeze since the FID is marked
18         a. n_bg_cnt <= bg_cnt
19         b. n_cell_cnt <= cell_cnt
20     3. else if(bg_cnt==0) // start a new counting cycle
21         a. n_bg_cnt <= bg_cnt_value
22         b. n_cell_cnt <= 0
23     4. else
24         a. n_bg_cnt <= bg_cnt - 1
25         b. n_cell_cnt <= cell_cnt
26     5. Write FID_mem[FID] <= {cell_cnt, mark, n_bg_cnt, bg_cnt_value}
27

```

Software read process

```

28
29     1. RD FID1_mem[MARK_RP] {cell_cnt, mark, bg_cnt, bg_cnt_value}
30     2. RD FID2_mem[MARK_RP] {next_FID}
31     3. n_cell_cnt = 0
32     4. n_mark = 0
33     5. n_bg_cnt = bg_cnt_value
34     6. if(MARK_EMPTY==0) // some marked FIDs are avail
35         if(MARK_RP==MARK_WP)// last FID of marked FID's link list
36             n_MARK_RP <= MARK_RP
37             n_MARK_EMPTY <= 1
38         else // more than one marked FID's avail
39             n_MARK_RP <= next_FID
40             n_MARK_EMPTY <= 0
41     else // no marked FID's avail
42         n_MARK_RP <= MARK_RP
43         n_MARK_EMPTY <= 1
44     7. if(VALID==1 & MARK_EMPTY==0)
45         WR FID1_mem[MARK_RP] <= {n_cell_cnt, n_mark, n_bg_cnt,
46         bg_cnt_value}

```

Calculation of threshold for Meter mode

The meter block can measure the number of cells for different intervals of times depending on the value of the BG_CNT[3:0] as shown in Figure 422,

For each measured output rate there is a maximum amount of cells that can be accumulated in 40mSec amount of time (threshold for 40mSec). The equation is:

$$\text{Threshold} [\# \text{Cells}/40\text{mSec}] = \text{Rout}[\text{bits}/\text{Sec}]/(25*512) = \text{Rout}/12800$$

For example, a rate of 64,000bps means a threshold of 5 cells/40mSec. Since this number is too small, it doesn't enable much of granularity in the measurement, an increase on BG_CNT can help since the measurement will be over longer period of time. If BG_CNT is 15, then the threshold should increase to $5*16 = 80$ cells.

The user has the tradeoff between the time to measure and the granularity of the measurement.

Memory description

Figure 423 shows what memories and bits are needed per mode of operation.

Rate ID/Threshold Internal Memory

Shaper Rate ID for Shaper mode

Length: 1024, Address: 10 bit rateID from database block, #bits: 86
See Figure 424.

Meter Threshold for Meter mode

Length: 1024, Address: 10 bit ThrsID from database block, #bits: 27
See Figure 425.

Shaper Slot Memory

Length: 256k, Address: 18 bit slot location, #bits: 41
See Figure 426.

Shaper/Meter FlowID Memory

The FID memory is divided to 2 memory blocks with different address to each. Word write-enable signals are assigned to these memories for the shaper mode. Meter algorithm uses 36 bits in FID1 memory and no bits in FID2 memory. Shaper algorithm uses 65 bits in FID1 memory and 25 bits in FID2 memory. When using single-leaky-bucket only, the unused memory bits should be tied down to zero.

FID MEM1 for Shaper

Length: 1M, Address: 20 bit FlowID, #bits 72

1 See Figure 427. (If more bits are needed because of future needs, the credit field
2 can be reduced)

3

4 **FID MEM1 for Meter**

5 Length: 1M, Address: 20 bit FlowID, #bits: 36

6 See Figure 428. There is no need to word write-enable signals in case of Meter
7 mode.

8

9 **FID MEM2 for Shaper**

10 Length: 1M, Address: 20 bit FlowID, #bits: 24

11 See Figure 429.

12

13 **FID MEM2 for Meter**

14 Length: 1M, Address: 20 bit FlowID, #bits: 20

15 See Figure 430.

16 **CPU Interface**

17

18 **Registers Description**

19 See Figure 431.

20

21 **CPU Commands for Shaper mode**

22

23 **Read FID1 memory**

24 See Figure 432.

25 P – peak_sustain

26 Registers R2, R3 are used in case of dual-leaky-bucket only.

27

28 **Write FID1 memory when using shaper**

29 See Figure 433.

30 P – peak_sustain

31 Registers R2, R3 are used in case of dual-leaky-bucket only.

1

2 **Read FID2 Memory**

3 See Figure 434.

4 SBT[0] – start (bit 20 in the register)

5 SBT[2:1] – bg_stt[1:0] (bits [22:21] in the register)

6 SBT[3] – time_msb (bit [23] in the register)

7

8 **Write FID2 Memory**

9 See Figure 435.

10 SBT[0] – start (bit 20 in the register)

11 SBT[2:1] – bg_stt[1:0] (bits [22:21] in the register)

12 SBT[3] – time_msb (bit [23] in the register)

13

14 **Read SHP_SLOT memory**

15 See Figure 436.

16 E – empty indication (bit [20] in the register)

17

18 **Write SHP_SLOT memory**

19 See Figure 437.

20 E – empty indication (bit [20] in the register)

21

22 **Read RATE_ID memory**

23 See Figure 438.

24 Ks[17:0] – Ks Mantissa, Ks[22:18] – Ks Exponent (same for Kp)

25

26 **Write RATE_ID memory**

27 See Figure 439.

28 Ks[17:0] – Ks Mantissa, Ks[22:18] – Ks Exponent (same for Kp)

1

2 **Setup connection command**

3 See Figure 440. This command resets all fields in FID1 and FID2 memories.
4 Except the back_count field in FID2 memory, which is driven from the command itself.
5 Also bg_stt[1:0] and empty bits should be set to 1 during this command.

6

7 **CPU Commands for Meter mode**

8

9 **Read FID1 memory**

10 See Figure 441.
11 M – Mark bit (bit [27] in the register)
12 Cur[3:0] – bg_cnt
13 Cnt[3:0] – bg_cnt_value

14

15 **Write FID1 memory**

16 See Figure 442.

17

18 **Read FID2 Memory**

19 See Figure 443.

20

21 **Write FID2 Memory**

22 See Figure 444.

23

24 **Read Threshold (internal) Memory**

25 See Figure 445.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36

Write Threshold (internal) Memory

See Figure 446.

ReadModWrite Meter

See Figure 447.

M – Mark bit (bit [27] in the register)

Cur[3:0] – bg_cnt

Cnt[3:0] – bg_cnt_value

This command triggers a write back to the same location clearing all values to zero (except the bg_cnt_value). That's is why this command has a different opcode then a simple read.

Setup connection command

See Figure 448.

This command writes the following data to FID1 memory:

cell_cnt <= 0, mark <= 0, bg_cnt <= cnt[3:0], bg_cnt_value <= cnt[3:0]

Reserved Op-codes

The op-codes 1000, 1001, 1010, 1011, 1100, 1101, 1110, 1111 are reserved.

Timing Diagrams for Shaper

The Shaper is able to process one input phase, one background phase and one output phase every 8 clock cycle slots.

Timing Diagrams for Meter

The Meter needs to perform one input phase, one background process and one CPU Rd/Wr accesses during 8 cycles budget. Figure 449 is a timing diagram that shows the accesses to the memories during these 8 cycles.

The CPU Rd/Wr operation and the Meter software phase cannot exist together at the same 8 cycles. If a ReadModWrite command is issued then no other CPU interface is possible during these 8 cycles.

For regular Read/Write from the internal/External memories, all accesses occur at slot 7. There are no data dependencies between any of the accesses. The ReadModWrite command clears all fields therefore the background process can be overwritten by it.

MX 2 CPU Interface

This block is responsible for three different functions:

- Connects between the blocks and the CPU via the CPU interface. It decodes the address to Figure which block to access and generate the chip select to each block.
- Generates global synchronization signals to the rest of the blocks.
- Holds general control registers like version register, mode register etc.

Features

- Address bus of 10 bits per chip.
Bits 0-5 are used as an address inside the block (up to 64 registers per block)
Bits 9-6 are used to select a block (up to 16 blocks per chip)
- Data bus of 32 bits per chip.
Externally this bus is b-directional
Internally this bus is split to 2 busses: DATA_IN and DATA_OUT.
- Chip select signal is the qualifier for data, address and read/write busses.
- The CPU interface can support multiplexed CPU interface or non-multiplexed CPU interface. When using multiplexed interface, the address and data are folded to the DATA[31:0] bus.
- Support little endian and big endian.
- CPU accesses all registers in all blocks via direct read/write commands.
- CPU accesses all internal/external memories in all blocks via indirect accesses using the general-purpose registers dedicated for these accesses.
- Blocks addresses are shown in Figure 450.
- Generates a global sync signal once every 520 cycles of the 200MHz clock (65 ports, 8 cycles per port).
- Multiplexes all test busses from the blocks to one output test bus.

Functional Description

The CPU IF block is used to interface with the CPU. The CPU IF block decodes the incoming address bus to generate a unique chip-select signal to each block. The CPU IF block also holds general registers that contain general parameters and modes for the entire chip. Basically any logic that is not related directly to a specific block is in this block.

The CPU IF block can operate in two modes depending on an input pins to select between them:

- Regular CPU access: separate busses for data and for address. In this case address bus is qualified by the chip-select signal.
- Multiplexed CPU access: one bus is multiplexed for address and for data. In this case address is qualified by the ALE signal.

The CPU interface is based on Intel's microprocessor interface (when multiplexing address/data bus). The following timing diagrams defines the characteristics of the CPU interface for Read and for Write accesses:

Read access characteristics

See Figure 451.

- 1 • RDWR_L signal has the same timings as the ADDR bus.
- 2 • A valid read cycle is defined when RDWR_L is latched as high on the negative edge
- 3 of the CS_L signal.
- 4 • In non-multiplexed address/data bus architecture, ALE should be held high so
- 5 parameters t_{Salr} , t_{Halr} , t_{VL} , t_{Slr} and t_{Hlr} are not applicable.
- 6 • In our implementation the address is latched on the negative edge of the CS_L
- 7 therefore t_{Har} is not applicable.
- 8 See Figure 452.

10 Write access characteristics

11 See Figure 453.

- 12 • RDWR_L signal has the same timings as the ADDR bus.
- 13 • A valid write cycle is defined when RDWR_L is latched as low on the negative edge
- 14 of the CS_L signal.
- 15 • In non-multiplexed address/data bus architecture, ALE should be held high so
- 16 parameters t_{Salw} , t_{Halw} , t_{VL} , t_{Slw} and t_{Hlw} are not applicable.
- 17 • In our implementation the address is latched on the negative edge of the CS_L
- 18 therefore t_{Haw} is not applicable.
- 19 See Figure 454.

21 Timing diagram of signals to blocks

22 The following diagram describes the timings of the CPU interface to each block.
 23 Note that chip-select signal is now asserted for only one cycle of the 200MHz clock. For
 24 a write or read cycle, all parameters are valid and stable before chip-select is asserted.
 25 When a read cycle occurs data is driven from the block after a maximum of 3 clock
 26 cycles of 200MHz clock.

27 See Figure 455.

29 Configuration modes

30 The CPU interface has 2 configuration bits to control the CPU interface as shown
 31 in Figure 456.

33 GSYNC - Synchronization signal

34 GSYNC signal is a one-cycle pulse in the 200MHz clock, which is asserted once
 35 every 520 cycles. The SYNC signal is used for synchronization between the blocks and
 36 between the chips as well as to generate the slot-count and the CPU-port indication.
 37 See Figure 457.

38 Slot count and CPU port indication

39 The signals slot_count and CPU_port are generated inside each block base on the
 40 GSYNC signal. Each block first needs to sample GSYNC and use the sampled version to
 41 generate the signals. See Figure 458.

1 **Test Multiplexing**

2 The CPU interface block receives test busses from each block and multiplex them
3 to one output test bus. The width of all test busses is 32 bits. Each block has a control
4 register to select the information to drive to the pins. The CPU interface block has a
5 control register to select between the blocks.
6 See Figure 459.

8 **Interface**

9 The CPU IF block generates chip-select signal per block, the rdwr_l, addr and
10 data busses are propagated to the blocks directly. The blocks perform the read or write
11 command based on the rdwr_l signal only when the chip-select signal per the block is
12 asserted. The chip-select signal per each block is synchronous to the 200MHz clock
13 domain already. Each block first samples the incoming signals and then uses them. This
14 is being done to solve any timing issues that may occur. If during placement and routing
15 the signals from the CPU IF to the blocks will violate the 5nSec budget, another set of
16 FFs can be placed in the middle of the path.
17 For a read cycle the CPU IF block selects the data busses coming in from the blocks
18 based on the chip-select per block (the address).

19 See Figure 460; software reset and testing signals are not included in the diagram.
20 The address, data and rd/wr signals takes about 4 cycles of 200MHz to propagate from
21 the chip's pins to the block (6 cycles are used for delay calculation). 2 cycles are needed
22 to perform read or write accesses in the block.
23 Therefore a total of 8 cycles are needed to perform a read or a write access.

25 **Block Descriptions**

27 **CPU Interface**

28 The negative edge of the CS_L signal qualifies the ADDR[9:0] and the RDWR_L
29 signals. The positive edge of the CS_L signal qualifies the DATA bus for read and for
30 write accesses. The CS_L is basically an enable signal to latches. Figure 461 shows the
31 basic latching circuit and the synchronization to 200MHz clock domain:

33 **Registers**

34 Figure 462 describes the registers implemented in the CPU IF block. These
35 registers act as control to the rest of the chip. The registers are implemented in the
36 200MHz clock domain even if there is no need to do so since their value is fixed and
37 stable during normal operation.

39 **Appendix A: Tags for Lookup Engine**

40 The following types of traffic are supported on the Maximus chip Lookup Engine
41 Block:

- 42 1. ATM
- 43 2. MPLS
 - 44 a. ATM
 - 45 b. PPP
 - 46 c. Ethernet

1 d. Frame Relay

2 3. Ethernet

3 4. IP Frame Relay

4 (All bit numbering assumes Big Endian ordering)

5 See Figure 268.

6 For the first version of the chip, the Flow Descriptor(FD) in Lookup Engine block
7 will support Data Types 1 and 2(a) thru 2(d). Pure Ethernet, IP and Frame Relay will be
8 supported in the next version of the Maximus chip. Header format references for these
9 traffic types are given below:

10

11 **HEADER FORMATS**

12 **1. ATM – Tag in the Link Layer (carried in 28 bit fields - 12 bit VPI/16 bit VCI)**
13 (Computer Networks, 3rd Edition, Andrew Tannenbaum, pg. 451)

14 See Figure 269. Each field is 8 bits wide with 12 bit VPI and 16 bit VCI field, for
15 a total of 28 bits lookup tag. This 28 bit lookup tag is located in the 1st 64-bit word, bits 1
16 thru 28.

17

18 **2(a) MPLS – ATM, 2(b) MPLS-PPP, 2(c) MPLS-Ethernet and 2(d) MPLS-FR**
19 (Tag is always between the Link and the Network Layer)

20 (MPLS, Technology & Applications, Davie & Rekhter, pg. 132)

21 The MPLS Header is inserted between the Link Layer and the Network Layer of
22 each of the above four protocols. The detailed Link Layer description of each of the
23 protocols is given below in other sections of this document.

24 The MPLS field is the same for the four different Link Layer protocols. The exact
25 location of this 20-bit MPLS label in the Frame varies with the Link Layer used, and is
26 summarized in the table above. With 20 bits, the MPLS label allows having 2^{20} possible
27 tags. See Figure 270.

28

29 **3. Ethernet – Tag in the Link Layer**

30 (Computer Networks, 3rd Edition, Andrew Tannenbaum, pg. 281)

31 Ethernet Destination Address(DA) field is used in the Lookup Engine Block.
32 These are the 48 bits after the first 64 bits as shown below. See Figure 271.

33

34 **4. IP – Tag in Network Layer**

35 (Computer Networks, 3rd Edition, Andrew Tannenbaum, pg. 413)

36 The IP Packet Destination Address is used for Lookup Engine. Starting from the
37 the 3rd 64 bit word, the Destination Address is 32 bits field and is used in the Lookup
38 Engine. This assumes Big Endian Transmission. See Figure 272.

39

40 **5. Frame Relay (FR) - Tag in Link Layer (carried in 10 bit DLCI field)**

41 (Frame Relay Technology & Practice, Buckwalter, pg. 39) See Figure 273.

42

43 **6. Point to Point Protocol (PPP)**

44 (Computer Networks, 3rd Edition, Andrew Tannenbaum, pg. 232) See Figure 274.

Appendix B: Flow Control Mechanism of CSIX-L1

CSIX provides multiple levels of flow control and supports flow control in both directions. At ingress, flow control is fabric-to-TM; at egress, flow control is TM-to-fabric.

Link-level Flow Control

The link level flow control is symmetric across the CSIX interface and provides independent control for data and control queues. There are 2 link level queues. One queue is for control traffic and the other queue is for data traffic. For each queue there is a **ready** bit in every header indicating the congestion status for the receive queue of the respective traffic type. See Figure 275.

Fabric Flow Control

The fabric flow control provides a finer level of flow control than link level flow control by specifying the specific TM and Class that are oversubscribed. Flow control on a TM port (line-end or VC) basis is the responsibility of cooperating TMs. Flow control frames can go in both the ingress and egress directions. See Figure 276 and 277.

Appendix C—MX2 Software Interface: A User Guide to Software Interface

This document describes the software interface to the MX2 chip via the CPU interface. This document should include all atomic CPU commands as well as the high level commands.

There are two types of commands:

1. Low-level commands, when issued, the driver generates a single command to access registers or internal/external memories in one specific block.
For example read scheduler's external memory
2. High-level commands, when issued, the driver generates several commands to access registers or internal/external memories in several blocks.
For example setup connection command to initialize FID database in several blocks.

Features

1. the interface supports up to 16 different blocks in the chip, each block has up to 64 registers.
2. CPU accesses all registers in all blocks via direct read/write commands.
3. CPU accesses all internal/external memories in all blocks via indirect accesses using the general-purpose registers dedicated for these accesses.
4. CPU interface can support multiplexed address/data bus or non-multiplexed.
5. CPU interface can support big endian or little endian modes.

Interface

Physical Interface

The CPU interface to the MX2 chip is via standard CPU interface that contains signals as shown in Figure 463. The software, using this interface, can directly access the registers inside each block and can indirectly access internal and external memories per each block.

Address spaces per block/chip

The 4 MSBs of the address bus points to the specific block, while the 6 LSBs of the address points to internal registers inside the block. Each block has address space of 64 entries only. The internal registers of each block are directly mapped to these addresses and external/internal memories are indirectly mapped. The MX2 chip can support up to 16 blocks with 64 registers each. Figure 464 is the allocation of address spaces per block.

Direct Access

A direct access is used to access registers only. The only two commands used in a direct access are read and write. Each register has it's own unique address inside the chip.

Indirect Access

An indirect access is used to access internal or external memories. Each block has a set of registers dedicated for the indirect access. The COM register contains an Op-code and address of the memory to access. Registers R0 through Rn are used to contain the data.

Indirect Write Command

A write command starts with writing the data to registers R0 to Rn and then writing the opcode and the address to the COM register. The opcode defines which memory to access and if it is a read or write command.

The hardware then, takes the data from registers R0 through Rn and write them to a specific address defined in the COM register in the specific memory defined by the opcode.

Once the command is complete, the hardware resets the opcode bits in the COM registers. It is guarantee that aread/write command to internal/external memory will be complete after 12 cycles of the operating clock (@200MHz it means 60nSec).

Indirect Read Command

A read command starts with writing the opcode and the address to the COM register. The opcode defines which memory to access and if it is a read or write command.

1 The hardware then, uses the address to perform a read access the memory. Once
2 the data arrives the hardware assign the data to registers R0 through Rn and then clears
3 the opcode bits in the COM register.

4 The software then, generates direct reads to registers R0 through Rn to get the
5 data. It is guarantee that aread/write command to internal/external memory will be
6 complete after 12 cycles of the operating clock (@200MHz it means 60nSec).

7 **CPUIF Block**

9 CPUIF block doesn't have any internal or external memories, therefore, only
10 direct commands to registers are possible in this block.

12 **Registers**

13 See Figure 465.

15 **Commands**

16 N/A

18 **PFQ Block**

20 **Registers**

21 See Figure 466.

23 **Commands**

24 **Read FID Enqueue Memory (72 bits) – Up to 2M locations**

25 See Figure 467.

26 **Write FID Enqueue Memory (72 bits) – Up to 2M locations**

27 See Figure 468.

28 **Read FID Dequeue Memory (36 bits)– Up to 1M locations**

29 See Figure 469.

30 **Write FID Dequeue Memory (36 bits) – Up to 1M locations**

31 See Figure 470.

32 **Read BID Memory (36 bits) – Up to 16M locations**

33 See Figure 471.

1 **Write BID Memory (36 bits) – Up to 16M locations**

2 See Figure 472.

3 **Read Stat Memory (72 bits) – Up to 2M locations**

4 See Figure 473.

5 **Write Stat Memory (72 bits) – Up to 2M locations**

6 See Figure 474.

7 **Setup connection command for FID**

8 See Figure 475. This command sets up a connection by accessing external
9 memories and initializes the corresponding locations to the proper parameters. It is a
10 write command. Clear PFQ entries for ENQ, DEQ and STAT!

11 **Tear-Down connection command for FID**

12 See Figure 476. This command tears down a connection by accessing external
13 memories. It is a write command. The data is all “0” PFQ returns all remaining BIDS on
14 the queue of the down connection to the free link list during this command.

15 **INIT FID MEMORIES**

16 The memories are the enqueue, dequeue, and statistics memories. See Figure 477.

17 **INIT BID Memory**

18 This will build the free buffer list. See Figure 478.

19 **Reserved Op-codes**

20 The op-codes 1101, 1110, and 1111 are reserved. The following commands are
21 combinations of the above simple commands. A user command like this should be broken
22 to several Read or Write commands in the driver.

23 **Initialized External memories**

- 24 • Write “0” in all memory locations.
25 • Write “NULL” in the enqueue memory.
26 • Write “NULL” in the dequeue memory.
27 • Build the free buffer list by linking each buffer to the next one.

1 **Initialized Internal memories**

- 2 • Write "NULL" in the Multicast Root memory.
- 3 • Write "NULL" in the Multicast Leaf memory.
- 4 • Write "NULL" in the Tunneling Root memory.
- 5 • Write "NULL" in the Tunneling Leaf memory.

7 **DataBase Block**

9 **Registers**

10 See Figure 479.

12 **Commands**

13 **Read external FID memory**

14 See Figure 480.

- 15 A[0] – C_P (bit [6] in the register)
- 16 A[1] – CPU_PORT (bit [7] in the register)
- 17 Pri[2:0] – Priority for shaped traffic
- 18 S – Shape (bit [23] in the register)
- 19 C[0] – clpi (bit [10] in the register)
- 20 C[1] – Empty (bit [11] in the register)

21 **Write external FID memory**

22 See Figure 481.

- 23 A[0] – C_P (bit [6] in the register)
- 24 A[1] – CPU_PORT (bit [6] in the register)
- 25 Pri[2:0] – Priority for shaped traffic
- 26 S – Shape (bit [23] in the register)
- 27 C[0] – clpi (bit [10] in the register)
- 28 C[1] – Empty (bit [11] in the register)

29 **Setup FID connection**

30 See Figure 482.

- 31 A[0] – C_P (bit [6] in the register)
- 32 A[1] – CPU_PORT (bit [6] in the register)
- 33 Pri[2:0] – Priority for shaped traffic
- 34 S – Shape (bit [23] in the register)
- 35 During the setup connection command, the user defines the above fields only.
- 36 The fields RP, WP, NCELL, CELL_CNT and CLP should be reset to 0 by the hardware
- 37 The field EMPTY should be set to 1 by the hardware.

1 Read external PKT memory

2 See Figure 483.

3 L[0] – CLP

4 L[1] - Reserved

5 Write external PKT memory

6 See Figure 484.

7 L[0] – CLP

8 L[1] - Reserved

9 Read port calendar memory

10 See Figure 485.

11 Data[5:0] – portid

12 Data[6] – valid

13 Data[7] – jump

14 Write port calendar memory

15 See Figure 486.

16 Data[5:0] – portid

17 Data[6] – valid

18 Data[7] – jump

19 Read SHP Strict memory

20 See Figure 487.

21 P[2:0] – previous QOS pointer

22 P[3] – previous QOS pointer valid

23 Write SHP Strict memory

24 See Figure 488.

25 Read SHP out memory

26 See Figure 489.

27 Write SHP out memory

28 See Figure 490.

1 Init SHP_out_port memory

2 See Figure 491. This command generates write accesses to all entries in the
3 Shp_out memory (64 entries) with a specific data of: p_rp = 0, p_wp = 0, p_e = 1

4 Init Shp_strict memory

5 See Figure 492. This command generates write accesses to all entries in the
6 Shp_strict memory (64 entries) with a specific data of: empty = 8'b1111_1111,
7 prev_qos_v = 0, prev_qos = 0

8 Reserved Opcodes

9 The opcodes: 1110, 1111 are reserved.

11 Shaper/Meter Block**13 Registers**

14 See Figure 493.

16 Commands for Shaper mode**17 Read FID1 memory**

18 See Figure 494.
19 P – peak_sustain
20 Registers R2, R3 are used in case of dual-leaky-bucket only.

21 Write FID1 memory when using shaper

22 See Figure 495.
23 P – peak_sustain
24 Registers R2, R3 are used in case of dual-leaky-bucket only.

25 Read FID2 Memory

26 See Figure 496.
27 SBT[0] – start (bit 20 in the register)
28 SBT[2:1] – bg_stt[1:0] (bits [22:21] in the register)
29 SBT[3] – time_msb (bit [23] in the register)

30 Write FID2 Memory

31 See Figure 497.
32 SBT[0] – start (bit 20 in the register)
33 SBT[2:1] – bg_stt[1:0] (bits [22:21] in the register)

- 1 SBT[3] – time_msb (bit [23] in the register)
- 2 **Read SHP_SLOT memory**
- 3 See Figure 498.
- 4 E – empty indication (bit [20] in the register)
- 5 **Write SHP_SLOT memory**
- 6 See Figure 499.
- 7 E – empty indication (bit [20] in the register)
- 8 **Read RATE_ID memory**
- 9 See Figure 500.
- 10 Ks[17:0] – Ks Mantissa, Ks[22:18] – Ks Exponent (same for Kp)
- 11 **Write RATE_ID memory**
- 12 See Figure 501.
- 13 Ks[17:0] – Ks Mantissa, Ks[22:18] – Ks Exponent (same for Kp)
- 14 **Setup connection command**
- 15 See Figure 502.
- 16 This command resets all fields in FID1 and FID2 memories. Except the back_count field
- 17 in FID2 memory, which is driven from the command itself.
- 18 Also bg_stt[1:0] and empty bits should be set to 1 during this command.
- 19
- 20 **Commands for Meter mode**
- 21 **Read FID1 memory**
- 22 See Figure 503.
- 23 M – Mark bit (bit [27] in the register)
- 24 Cur[3:0] – bg_cnt
- 25 Cnt[3:0] – bg_cnt_value
- 26 **Write FID1 memory**
- 27 See Figure 504.
- 28 **Read FID2 Memory**
- 29 See Figure 505.

1 Write FID2 Memory

2 See Figure 506.

3 Read Threshold (internal) Memory

4 See Figure 507.

5 Write Threshold (internal) Memory

6 See Figure 508.

7 ReadModWrite Meter

8 See Figure 509.

9 M – Mark bit (bit [27] in the register)

10 Cur[3:0] – bg_cnt

11 Cnt[3:0] – bg_cnt_value

12 This command triggers a write back to the same location clearing all values to
13 zero (except the bg_cnt_value). That's is why this command has a different opcode then a
14 simple read.

15 Setup connection command

16 See Figure 510. This command writes the following data to FID1 memory:

17 cell_cnt <= 0, mark <= 0, bg_cnt <= cnt[3:0], bg_cnt_value <= cnt[3:0]

18 Reserved Op-codes

19 The op-codes 1000, 1001, 1010, 1011, 1100, 1101, 1110, 1111 are reserved.

20
21 Scheduler Block**22**
23 Registers

24 See Figure 511.

25 Commands**26 Read FID NEXT Memory**

27 See Figure 512.

28 Write FID NEXT Memory

29 See Figure 513.

1 Read QOS Parameters

2 See Figure 514.

3 Write QOS Parameters

4 See Figure 515.

5 Read QOS Descriptor

6 See Figure 516.

7 Write QOS Descriptor

8 See Figure 517.

9 Read PORT Parameters

10 See Figure 518.

11 Write PORT Parameters

12 See Figure 519.

13 Read PORT Descriptor

14 See Figure 520.

15 PQ[0] = PREV_QOS_VALID

16 PQ[3:1] = PREV_QOS[2:0]

17 PQ[6:4] = ACTIVE_PTR[2:0]

18

19 Write PORT Descriptor

20 See Figure 521.

21 Read PORT Next

22 See Figure 522.

23 Write PORT Next

24 See Figure 523.

1 Init QOS

2 See Figure 524. This command write the parameters to the QOS Parameter
3 Memory and also initiates a write to the QOS Descriptor Memory with the following
4 values:

5 F_EMPTY = 1

6 Rest of the fields should be reset to zero

7 Init port

8 See Figure 525. This command write the parameters to the Port Parameter
9 Memory and also initiates a write to the Port Descriptor Memory with the following
10 values:

11 QW_EMPTY[7:0] and QA_EMPTY[7:0] = 8'b1111_1111

12 Rest of the fields are reset to zero.

13 Reserved Op-codesReserved Op-codesReserved Op-codes

14 The op-codes 1111 are reserved.

16 Driver's high level commands

17 The following section describes higher level commands then just read/write to
18 internal/external memories. The following commands are basically combinations of the
19 above commands and are created to give easier user interface to the chip.

20 These commands doesn't have special opcode in the chip itself, they are implemented in
21 the driver only.

23 Setup Connection Command

24 This command sets up a connection by accessing external memories in several
25 blocks and initializes the proper parameters per FlowID.

26 The following table describes the parameters per FlowID that should be set during
27 the setup connection command, for more specific information about the fields, please
28 refer to the blocks SPECs in Figure 526.

29 Once the driver receives a setup connection command, the driver generates setup
30 connection command to the DBS, SHP and PFQ blocks. The parameters to the setup
31 connection commands to the blocks are driven from the user along with the global setup
32 connection command to the driver.

34 Tear-down Connection Command

35 A tear-down-connection command is used only in the PFQ block. This command
36 is needed to invalidate some fields in the FID descriptor memories as well as the statistics
37 fields. The Database, Shaper and Scheduler blocks don't need to invalidate any of the
38 FID fields. Therefore, there is no special driver command for this command. The user can
39 initiate a simple write command to the PFQ to implement this command.

Power-up initialization

When reset is applied to the chip (or power-up) the internal control registers wakes up in their default value, all traffic is disabled and internal/external memories are not defined. The following is the initialization sequence the user should perform before allowing any traffic to the MX2 chip:

1. PFQ – issue “Init FID Memory” command
2. PFQ – issue “Init BID Memory” command
3. PFQ – initialize multicast internal memories (if needed)
4. PFQ – initialize tunneling internal memory (if needed)
5. PFQ – specify bits in control register if different then default (address 36)
6. PFQ – initialize RED internal memory only if using RED
7. PFQ – initialize class threshold registers (addresses 45 to 52) only when using CLASS
8. DBS – initialize Port Calendar memory
9. DBS – issue “Init Shp_strict Memory” command
10. DBS – issue “Init Shp_out_port Memory” command
11. CPUIF – define Shaper/Meter mode (register 1 bit [8])
12. SHP – initialize RateID/Threshold memory
13. SHP – issue “Init Slot Memory” command
14. SCH – initialize Port Descriptor Memory
15. SCH – initialize QOS descriptor Memory
16. <setup connections if needed>
17. PFQ – enable input traffic for specific ports by writing to registers 43 and 44
18. SCH – enable output traffic from scheduler (register 32 bit [0])
19. SHP – enable output traffic from shaper (register 32 bit [1])

Once these steps are complete, the user can initiate setup connection commands.

The user doesn’t have to go through all of the above steps, for example, if tunneling is not used then there is no need to initialize the internal memory for tunneling in the PFQ block.the MX2 Chip

Appendix D—MX2 External Memory

See Figure 527.

While this invention has been described in terms of several preferred embodiments, there are alterations, permutations, and equivalents, which fall within the scope of this invention. It should also be noted that there are many alternative ways of implementing the methods and apparatuses of the present invention. For example, the present invention may be practiced by integrating line card functionality with switch fabric functionality such that a router does not have individual line cards. Alternatively, router in accordance with the invention may include many lines cards deploying the Maximus device, some of the line cards functioning as ingress line cards, others of the line cards functioning as egress line cards. It is therefore intended that the following appended claims be

- 1 interpreted as including all such alteration, permutations, and equivalents as fall within
- 2 the true spirit and scope of the present invention.

the present invention, the scope of the present invention is not limited by the above description.